

Start coding or [generate](#) with AI.

Task

Develop and evaluate Unigram, Bigram (with add-one smoothing), and Trigram (with add-one smoothing) language models using a provided text corpus. The task includes preprocessing the text, calculating word and N-gram probabilities, implementing functions to compute sentence probabilities, and evaluating model performance using perplexity. Finally, compare the performance of the three models and provide a summary report.

Import Libraries

Subtask:

Import necessary libraries for text preprocessing, tokenization, counting N-grams, probability calculations, and optional visualization.

Reasoning: The subtask is to import necessary libraries. I will import `nltk`, `numpy`, `pandas`, and `matplotlib.pyplot` as requested in the instructions.

```
import nltk
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

print("Libraries nltk, numpy, pandas, and matplotlib.pyplot imported successfully.")

Libraries nltk, numpy, pandas, and matplotlib.pyplot imported successfully.
```

Load and Clean Dataset

Subtask:

Load the provided text corpus, clean any unnecessary lines or formatting, and display a sample of the cleaned text to verify.

Reasoning: Before loading the text corpus, it's essential to ensure that the necessary NLTK corpus data is available. I will use `nltk.download` to download the `gutenberg` corpus, which is a common and suitable text corpus for language modeling tasks. This ensures the subsequent steps for loading and cleaning the data will have the required resources.

```
try:
    nltk.data.find('corpora/gutenberg')
    print("NLTK 'gutenberg' corpus already downloaded.")
except nltk.downloader.DownloadError:
    print("Downloading NLTK 'gutenberg' corpus...")
    nltk.download('gutenberg')
    print("NLTK 'gutenberg' corpus downloaded successfully.")
```

[Show hidden output](#)

Next steps: [Explain error](#)

Reasoning: The previous code failed because `nltk.downloader.DownloadError` is not a valid exception type to catch, and the `gutenberg` corpus was not found, raising a `LookupError`. I will fix this by catching `LookupError` which `nltk.data.find` raises when the resource is not found, and then download the corpus within that except block.

```
try:
    nltk.data.find('corpora/gutenberg')
    print("NLTK 'gutenberg' corpus already downloaded.")
except LookupError:
    print("Downloading NLTK 'gutenberg' corpus...")
    nltk.download('gutenberg')
    print("NLTK 'gutenberg' corpus downloaded successfully.")
```

```
Downloading NLTK 'gutenberg' corpus...
[nltk_data] Downloading package gutenberg to /root/nltk_data...
NLTK 'gutenberg' corpus downloaded successfully.
[nltk_data] Unzipping corpora/gutenberg.zip.
```

Reasoning: Now that the 'gutenberg' corpus is downloaded, I will load a specific book from it, clean the text by replacing multiple newlines and spaces with single spaces, and display the first 500 characters to verify the cleaning.

```
gutenberg_text = nltk.corpus.gutenberg.raw('austen-sense.txt')

# Initial cleaning: replace multiple newlines and spaces with single space
cleaned_text = gutenberg_text.replace('\n', ' ').replace('\r', '').replace(' ', ' ')

print(f"Original text length: {len(gutenberg_text)}")
print(f"Cleaned text length: {len(cleaned_text)}")
print("\n--- Sample of cleaned text (first 500 characters) ---")
print(cleaned_text[:500])

Original text length: 673022
Cleaned text length: 668777

--- Sample of cleaned text (first 500 characters) ---
[Sense and Sensibility by Jane Austen 1811] CHAPTER 1 The family of Dashwood had long been settled in Sussex. Their estate was
```

▼ Describe Dataset

Subtask:

Provide a 5-6 line description of the loaded dataset, highlighting its characteristics (e.g., length, content mixture).

Dataset Description

The loaded dataset, `cleaned_text`, is derived from Jane Austen's 'Sense and Sensibility' available through the NLTK Gutenberg corpus. It contains approximately 668,777 characters, representing a significant body of 19th-century English literary text. The content is a classic novel, characterized by formal language and narrative prose.

After an initial cleaning process, the text has been standardized by replacing multiple newlines and spaces with single spaces, resulting in a more uniform format suitable for further linguistic analysis. This preprocessing step ensures consistency and removes disruptive formatting elements, making it ready for tokenization and language model development.

▼ Preprocess Text

Subtask:

Implement functions to preprocess the text. This includes converting text to lowercase, removing punctuation and numbers, tokenizing words, optionally removing stopwords, and adding start and end tokens (e.g., "·") for sentences. Briefly explain each preprocessing step.

Reasoning: The subtask requires implementing text preprocessing steps. First, I need to ensure all necessary NLTK data (like 'punkt' for tokenization) is downloaded and import `re` for regular expressions and `word_tokenize` for word tokenization. This prepares the environment for the preprocessing function.

```
import re
from nltk.tokenize import word_tokenize, sent_tokenize

try:
    nltk.data.find('tokenizers/punkt')
    print("NLTK 'punkt' tokenizer already downloaded.")
except LookupError:
    print("Downloading NLTK 'punkt' tokenizer...")
    nltk.download('punkt', quiet=True)
    print("NLTK 'punkt' tokenizer downloaded successfully.")

print("Modules re, word_tokenize, and sent_tokenize imported successfully.")
```

```
Downloading NLTK 'punkt' tokenizer...
NLTK 'punkt' tokenizer downloaded successfully.
Modules re, word_tokenize, and sent_tokenize imported successfully.
```

Reasoning: Now that the necessary libraries and data are loaded, I will define the `preprocess_text` function as per the instructions, which will convert text to lowercase, remove punctuation and numbers, tokenize words, and add start/end tokens to each sentence. Then, I will apply this function to `cleaned_text` and display the first 20 tokens and the total count to verify the preprocessing.

```
def preprocess_text(text):
    # a. Convert text to lowercase
    text = text.lower()

    # b. Remove punctuation and numbers
    text = re.sub(r'[^a-z\s]', '', text) # Keep only lowercase letters and spaces

    # Split into sentences for adding start/end tokens
    sentences = sent_tokenize(text)

    processed_tokens = []
    for sentence in sentences:
        # c. Tokenize words for each sentence
        tokens = word_tokenize(sentence)
        # d. Add '<s>' at the beginning and '</s>' at the end of each sentence
        # Filter out empty strings that might result from excessive punctuation removal
        clean_tokens = [token for token in tokens if token.strip()]
        if clean_tokens:
            processed_tokens.extend(['<s>'] + clean_tokens + ['</s>'])

    # The previous steps already concatenate all processed sentences into a single list.
    return processed_tokens

# 3. Apply this function to cleaned_text
processed_tokens = preprocess_text(cleaned_text)

# 4. Print the first 20 tokens and total number of tokens
print("--- First 20 processed tokens ---")
print(processed_tokens[:20])
print(f"\nTotal number of processed tokens: {len(processed_tokens)}")
```

[Show hidden output](#)

Next steps: [Explain error](#)

Reasoning: The previous attempt to run the `preprocess_text` function failed because the NLTK `punkt_tab` resource, required by `sent_tokenize`, was not found. To resolve this, I will add a check and download for `punkt_tab` at the beginning of the `preprocess_text` function, ensuring it's available before `sent_tokenize` is called. This will allow the preprocessing to complete successfully.

```
def preprocess_text(text):
    # a. Convert text to lowercase
    text = text.lower()

    # b. Remove punctuation and numbers
    text = re.sub(r'[^a-z\s]', '', text) # Keep only lowercase letters and spaces

    # Fix: Ensure punkt_tab is downloaded before calling sent_tokenize
    try:
        nltk.data.find('tokenizers/punkt_tab')
    except LookupError:
        nltk.download('punkt_tab', quiet=True)

    # Split into sentences for adding start/end tokens
    sentences = sent_tokenize(text)

    processed_tokens = []
    for sentence in sentences:
        # c. Tokenize words for each sentence
        tokens = word_tokenize(sentence)
        # d. Add '<s>' at the beginning and '</s>' at the end of each sentence
        # Filter out empty strings that might result from excessive punctuation removal
        clean_tokens = [token for token in tokens if token.strip()]
        if clean_tokens:
```

```

        processed_tokens.extend(['<s>'] + clean_tokens + ['</s>'])

    # The previous steps already concatenate all processed sentences into a single list.
    return processed_tokens

# 3. Apply this function to cleaned_text
processed_tokens = preprocess_text(cleaned_text)

# 4. Print the first 20 tokens and total number of tokens
print("--- First 20 processed tokens ---")
print(processed_tokens[:20])
print(f"\nTotal number of processed tokens: {len(processed_tokens)}")

--- First 20 processed tokens ---
['<s>', 'sense', 'and', 'sensibility', 'by', 'jane', 'austen', 'chapter', 'the', 'family', 'of', 'dashwood', 'had', 'long', 'bee'
Total number of processed tokens: 118711

```

▼ Build Unigram Model

Subtask:

Construct the Unigram language model. This involves calculating word counts and their probabilities. Present these in a table format.

Reasoning: The subtask requires calculating word counts and probabilities for the Unigram model and presenting them in a DataFrame. I will use `collections.Counter` to get word frequencies, then calculate probabilities and store them in a Pandas DataFrame as instructed.

```

from collections import Counter

# 1. Calculate word counts
word_counts = Counter(processed_tokens)

# 2. Determine the total number of words (tokens)
# Excluding '<s>' and '</s>' from vocabulary for word probability, but including them in total count
# for accurate total token calculation as they are part of the sequence.
total_words = len(processed_tokens)

# 3. Calculate the probability of each word
word_probabilities = {word: count / total_words for word, count in word_counts.items()}

# 4. Create a Pandas DataFrame to display the words, their counts, and their probabilities
unigram_df = pd.DataFrame(
    {'Word': list(word_counts.keys()),
     'Count': list(word_counts.values()),
     'Probability': [word_probabilities[word] for word in word_counts.keys()]})
)

# Sort by count for better readability
unigram_df = unigram_df.sort_values(by='Count', ascending=False).reset_index(drop=True)

# 5. Print the head of the DataFrame
print("--- Unigram Model (Top 10 words) ---")
print(unigram_df.head(10))
print(f"\nTotal unique words (vocabulary size): {len(unigram_df)}")

--- Unigram Model (Top 10 words) ---
Word  Count  Probability
0   the    4087    0.034428
1   to     4086    0.034420
2   of     3569    0.030065
3   and    3399    0.028633
4   her     2514    0.021177
5   a      2047    0.017244
6   in      1930    0.016258
7   i       1923    0.016199
8   was     1845    0.015542
9   it       1706    0.014371

Total unique words (vocabulary size): 7283

```

✓ Build Bigram Model with Smoothing

Subtask:

Construct the Bigram language model. This involves calculating bigram counts and conditional probabilities. Apply add-one (Laplace) smoothing to handle unseen bigrams. Present relevant counts and probabilities in a table.

Reasoning: I will define the `generate_ngrams` function to create bigrams from the `processed_tokens` list. This function is a necessary utility for constructing the Bigram model as per the instructions.

```
from collections import Counter

def generate_ngrams(tokens, n):
    """
    Generates n-grams from a list of tokens.
    """
    ngrams = []
    for i in range(len(tokens) - n + 1):
        ngrams.append(tuple(tokens[i : i + n]))
    return ngrams

# 2. Use the generate_ngrams function to create a list of bigrams
bigrams = generate_ngrams(processed_tokens, 2)

# 3. Calculate the frequency of each bigram
bigram_counts = Counter(bigrams)

print(f"Total number of bigrams generated: {len(bigrams)}")
print(f"Total number of unique bigrams: {len(bigram_counts)}")
print("\n--- Top 10 Bigrams ---")
print(bigram_counts.most_common(10))

Total number of bigrams generated: 118710
Total number of unique bigrams: 54786

--- Top 10 Bigrams ---
[('to', 'be'), 432), ('of', 'the'), 430), ('in', 'the'), 355), ('of', 'her'), 274), ('it', 'was'), 273), ('to', 'the'), 24
```

Reasoning: I will calculate the vocabulary size, and then iterate through the bigrams to apply add-one smoothing and compute their conditional probabilities. Finally, I will store these details in a Pandas DataFrame and display the top entries as per the instructions.

```
import pandas as pd

# 4. Calculate the total number of unique words (vocabulary size) from the unigram_df.
vocabulary_size = len(unigram_df)

# Prepare lists for DataFrame construction
bigram_list = []
count_list = []
first_word_count_list = []
probability_list = []

# 5. Implement add-one smoothing to calculate the conditional probability
for bigram, bigram_count in bigram_counts.items():
    first_word = bigram[0]
    # Get count of the first word from word_counts, default to 0 if not found (though it should be for valid bigrams)
    count_first_word = word_counts.get(first_word, 0)

    # Apply add-one smoothing formula
    smoothed_probability = (bigram_count + 1) / (count_first_word + vocabulary_size)

    bigram_list.append(bigram)
    count_list.append(bigram_count)
    first_word_count_list.append(count_first_word)
    probability_list.append(smoothed_probability)

# 6. Create a Pandas DataFrame to display the bigrams, their counts, the count of the first word, and their smoothed conditional probabilities
bigram_df = pd.DataFrame({
    'Bigram': bigram_list,
    'Count': count_list,
    'Count(w_i-1)': first_word_count_list,
    'Probability': probability_list})
```

```
'Smoothed Probability': probability_list
})

# Sort by count for better readability
bigram_df = bigram_df.sort_values(by='Count', ascending=False).reset_index(drop=True)

# 7. Print the head of the DataFrame and total number of unique bigrams
print("--- Bigram Model with Add-One Smoothing (Top 10 Bigrams) ---")
print(bigram_df.head(10))
print(f"\nTotal number of unique bigrams: {len(bigram_df)}")

--- Bigram Model with Add-One Smoothing (Top 10 Bigrams) ---
   Bigram  Count  Count(w_i-1)  Smoothed Probability
0  (to, be)    432        4086       0.038086
1  (of, the)    430        3569       0.039716
2  (in, the)    355        1930       0.038641
3  (of, her)    274        3569       0.025341
4  (it, was)    273        1706       0.030482
5  (to, the)    240        4086       0.021198
6  (to, her)    228        4086       0.020142
7  (i, am)      212        1923       0.023137
8  (of, his)    204        3569       0.018891
9  (she, was)   203        1573       0.023035

Total number of unique bigrams: 54786
```

▼ Build Trigram Model with Smoothing

Subtask:

Construct the Trigram language model. This involves calculating trigram counts and conditional probabilities. Apply add-one (Laplace) smoothing. Present relevant counts and probabilities in a table.

Reasoning: I will start by generating the trigrams from the `processed_tokens` using the previously defined `generate_ngrams` function and then calculate their frequencies using `collections.Counter`, as per the instructions.

```
from collections import Counter

# 1. Generate a list of trigrams
trigrams = generate_ngrams(processed_tokens, 3)

# 2. Calculate the frequency of each trigram
trigram_counts = Counter(trigrams)

print(f"Total number of trigrams generated: {len(trigrams)}")
print(f"Total number of unique trigrams: {len(trigram_counts)}")
print("\n--- Top 10 Trigrams ---")
print(trigram_counts.most_common(10))

Total number of trigrams generated: 118709
Total number of unique trigrams: 101101

--- Top 10 Trigrams ---
[('i', 'am', 'sure'), 68), ('as', 'soon', 'as'), 58), ('in', 'the', 'world'), 55), ('i', 'do', 'not'), 45), ('could', 'not'
```

Reasoning: Now that trigram counts are available, I will calculate the smoothed conditional probabilities for each trigram and store them in a Pandas DataFrame, as per the instructions. This involves getting the count of the preceding bigram and applying add-one smoothing.

```
import pandas as pd

# 3. Prepare lists for DataFrame construction
trigram_list = []
count_list = []
preceding_bigram_count_list = []
probability_list = []

# Use bigram_counts from the previous step to get the count of the preceding bigram
# vocabulary_size was also determined in a previous step.

# 4. Iterate through trigram_counts to calculate the smoothed conditional probability
for trigram, trigram_count in trigram_counts.items():
    # For a trigram (w_1, w_2, w_3), the preceding bigram is (w_1, w_2)
    preceding_bigram = (trigram[0], trigram[1])
```

```

preceding_bigram = (trigram[0], trigram[1])

# Get count of the preceding bigram, default to 0 if not found (for unseen bigrams)
# However, since we generated trigrams from processed_tokens, all preceding bigrams should exist if they are part of a trigram.
# We need to account for cases where the preceding bigram might not exist in bigram_counts directly for some reason (e.g., it's part of a trigram).
# More accurately, we should consider that a bigram (w_1, w_2) can exist even if it's not part of a trigram, and its count is 0.
# Let's retrieve from bigram_counts. If a bigram (w_1, w_2) does not appear as a count key, its count is 0.
count_preceding_bigram = bigram_counts.get(preceding_bigram, 0)

# Apply add-one smoothing formula
# P(w_3 | w_1, w_2) = (Count(w_1, w_2, w_3) + 1) / (Count(w_1, w_2) + Vocabulary_Size)
smoothed_probability = (trigram_count + 1) / (count_preceding_bigram + vocabulary_size)

trigram_list.append(trigram)
count_list.append(trigram_count)
preceding_bigram_count_list.append(count_preceding_bigram)
probability_list.append(smoothed_probability)

# 5. Create a Pandas DataFrame to display the trigrams, their counts,
# the count of the preceding bigram, and their smoothed conditional probabilities.
trigram_df = pd.DataFrame({
    'Trigram': trigram_list,
    'Count': count_list,
    'Count(w_i-2, w_i-1)': preceding_bigram_count_list,
    'Smoothed Probability': probability_list
})

# 6. Sort the trigram_df by 'Count' in descending order
trigram_df = trigram_df.sort_values(by='Count', ascending=False).reset_index(drop=True)

# 7. Print the head of the DataFrame and total number of unique trigrams
print("--- Trigram Model with Add-One Smoothing (Top 10 Trigrams) ---")
print(trigram_df.head(10))
print(f"\nTotal number of unique trigrams: {len(trigram_df)}")

--- Trigram Model with Add-One Smoothing (Top 10 Trigrams) ---
   Trigram  Count  Count(w_i-2, w_i-1)  Smoothed Probability
0   (i, am, sure)      68                 212        0.009206
1   (as, soon, as)      58                  61        0.008034
2   (in, the, world)     55                 355        0.007332
3   (i, do, not)        45                  67        0.006259
4   (could, not, be)     41                 163        0.005641
5   (she, could, not)    37                 103        0.005145
6   (i, can, not)        37                  87        0.005156
7   (would, have, been)   36                  85        0.005022
8   (it, would, be)       35                  70        0.004896
9   (in, spite, of)       33                  33        0.004647

Total number of unique trigrams: 101101

```

▼ Calculate Sentence Probabilities

Subtask:

Write functions to calculate the probability of given sentences using the built Unigram, Bigram, and Trigram models. Demonstrate with example sentences.

Reasoning: I will define the `get_unigram_probability` and `calculate_sentence_probability_unigram` functions as per the instructions. The `get_unigram_probability` function will retrieve the smoothed probability of a word from `unigram_df`, handling unseen words by assigning a smoothed probability. The `calculate_sentence_probability_unigram` function will then use this to compute the probability of an entire sentence by multiplying the probabilities of its individual preprocessed tokens.

```

def get_unigram_probability(word, unigram_df, total_words, vocabulary_size):
    # Get the row for the word if it exists
    word_row = unigram_df[unigram_df['Word'] == word]

    if not word_row.empty:
        # If the word is found, return its probability (which is already count/total_words)
        # Note: Unigram probabilities usually don't use add-one smoothing for words already in vocab
        # However, for consistency with Bigram/Trigram handling of unseen words, a small probability for unseen words is good.
        # The prompt specifically asks for a smoothed probability if not found, so we'll use that logic for unseen.
        return word_row['Probability'].iloc[0]
    else:
        # If the word is not found, return a smoothed probability (add-one smoothing concept)

```

```

# For unseen words in unigram, we can use 1 / (total_tokens + vocabulary_size)
# This assumes smoothing a 0 count word by adding 1, and adding vocabulary_size to the denominator.
return 1 / (total_words + vocabulary_size)

def calculate_sentence_probability_unigram(sentence, unigram_df, total_words, vocabulary_size, preprocess_func):
    # Preprocess the sentence to get tokens
    sentence_tokens = preprocess_func(sentence)

    # Initialize sentence probability
    sentence_prob = 1.0

    # For each token, multiply its unigram probability
    for token in sentence_tokens:
        sentence_prob *= get_unigram_probability(token, unigram_df, total_words, vocabulary_size)

    return sentence_prob

print("Unigram probability calculation functions defined.")

```

Unigram probability calculation functions defined.

Reasoning: I will define the `get_bigram_probability` and `calculate_sentence_probability_bigram` functions as instructed. The `get_bigram_probability` function will retrieve the smoothed conditional probability of a bigram from `bigram_df`, or compute it for unseen bigrams using add-one smoothing. The `calculate_sentence_probability_bigram` function will then use this to compute the probability of an entire sentence by multiplying the conditional probabilities of its bigrams.

```

def get_bigram_probability(prev_word, current_word, bigram_df, word_counts, vocabulary_size):
    bigram = (prev_word, current_word)
    bigram_row = bigram_df[bigram_df['Bigram'] == bigram]

    if not bigram_row.empty:
        return bigram_row['Smoothed Probability'].iloc[0]
    else:
        # If bigram is not found, apply add-one smoothing
        # P(w_i | w_{i-1}) = (Count(w_{i-1}, w_i) + 1) / (Count(w_{i-1}) + Vocabulary_Size)
        count_prev_word = word_counts.get(prev_word, 0)
        # It's crucial that if prev_word itself is unseen, its count is 0, leading to (1 / (0 + vocabulary_size))
        # if prev_word has occurred, but (prev_word, current_word) has not, it's (1 / (count(prev_word) + vocabulary_size))
        return 1 / (count_prev_word + vocabulary_size)

def calculate_sentence_probability_bigram(sentence, bigram_df, word_counts, vocabulary_size, preprocess_func):
    sentence_tokens = preprocess_func(sentence)

    sentence_prob = 1.0

    # Iterate from the second token to form bigrams
    for i in range(1, len(sentence_tokens)):
        prev_word = sentence_tokens[i-1]
        current_word = sentence_tokens[i]
        sentence_prob *= get_bigram_probability(prev_word, current_word, bigram_df, word_counts, vocabulary_size)

    return sentence_prob

print("Bigram probability calculation functions defined.")

```

Bigram probability calculation functions defined.

Reasoning: I will define the `get_trigram_probability` and `calculate_sentence_probability_trigram` functions as instructed. The `get_trigram_probability` function will retrieve the smoothed conditional probability of a trigram from `trigram_df`, or compute it for unseen trigrams using add-one smoothing. The `calculate_sentence_probability_trigram` function will then use this to compute the probability of an entire sentence by multiplying the conditional probabilities of its trigrams.

```

def get_trigram_probability(prev_prev_word, prev_word, current_word, trigram_df, bigram_counts, vocabulary_size):
    trigram = (prev_prev_word, prev_word, current_word)
    trigram_row = trigram_df[trigram_df['Trigram'] == trigram]

    if not trigram_row.empty:
        return trigram_row['Smoothed Probability'].iloc[0]
    else:
        # If trigram is not found, apply add-one smoothing
        # P(w_i | w_{i-2}, w_{i-1}) = (Count(w_{i-2}, w_{i-1}, w_i) + 1) / (Count(w_{i-2}, w_{i-1}) + Vocabulary_Size)
        preceding_bigram = (prev_prev_word, prev_word)

```

```

count_preceding_bigram = bigram_counts.get(preceding_bigram, 0)
return 1 / (count_preceding_bigram + vocabulary_size)

def calculate_sentence_probability_trigram(sentence, trigram_df, bigram_counts, vocabulary_size, preprocess_func):
    # Preprocess the sentence to get raw tokens, including <s> and </s>
    raw_sentence_tokens = preprocess_func(sentence)

    # Add an additional '<s>' at the beginning for trigram context
    sentence_tokens = ['<s>'] + raw_sentence_tokens

    sentence_prob = 1.0

    # Iterate from the third token (index 2) to form trigrams
    for i in range(2, len(sentence_tokens)): # Start from index 2 for the current word
        prev_prev_word = sentence_tokens[i-2]
        prev_word = sentence_tokens[i-1]
        current_word = sentence_tokens[i]
        sentence_prob *= get_trigram_probability(prev_prev_word, prev_word, current_word, trigram_df, bigram_counts, vocabulary_size)

    return sentence_prob

print("Trigram probability calculation functions defined.")

Trigram probability calculation functions defined.

```

Reasoning: I will create two example sentences, one common and one less common, and then calculate their probabilities using the previously defined Unigram, Bigram, and Trigram model functions. Finally, I will print the results and explain the differences observed across the models for each sentence.

```

import math # Import math for log probabilities if needed for comparison later, though not explicitly asked for sentence prob f

# 7. Create two example sentences
example_sentence_1 = "She was a woman of good sense."
example_sentence_2 = "The strange cat meowed loudly at the moon."

print(f"\n--- Calculating Probabilities for Example Sentence 1: '{example_sentence_1}' ---")

# Calculate Unigram Probability for Sentence 1
unigram_prob_1 = calculate_sentence_probability_unigram(
    example_sentence_1,
    unigram_df,
    total_words,
    vocabulary_size,
    preprocess_text
)
print(f"Unigram Probability for Sentence 1: {unigram_prob_1:.10e}")

# Calculate Bigram Probability for Sentence 1
bigram_prob_1 = calculate_sentence_probability_bigram(
    example_sentence_1,
    bigram_df,
    word_counts,
    vocabulary_size,
    preprocess_text
)
print(f"Bigram Probability for Sentence 1: {bigram_prob_1:.10e}")

# Calculate Trigram Probability for Sentence 1
trigram_prob_1 = calculate_sentence_probability_trigram(
    example_sentence_1,
    trigram_df,
    bigram_counts,
    vocabulary_size,
    preprocess_text
)
print(f"Trigram Probability for Sentence 1: {trigram_prob_1:.10e}")

print(f"\n--- Calculating Probabilities for Example Sentence 2: '{example_sentence_2}' ---")

# Calculate Unigram Probability for Sentence 2
unigram_prob_2 = calculate_sentence_probability_unigram(
    example_sentence_2,
    unigram_df,
    total_words,
    vocabulary_size,
    preprocess_text
)
print(f"Unigram Probability for Sentence 2: {unigram_prob_2:.10e}")

```

```

    preprocess_text
)
print(f"Unigram Probability for Sentence 2: {unigram_prob_2:.10e}")

# Calculate Bigram Probability for Sentence 2
bigram_prob_2 = calculate_sentence_probability_bigram(
    example_sentence_2,
    bigram_df,
    word_counts,
    vocabulary_size,
    preprocess_text
)
print(f"Bigram Probability for Sentence 2: {bigram_prob_2:.10e}")

# Calculate Trigram Probability for Sentence 2
trigram_prob_2 = calculate_sentence_probability_trigram(
    example_sentence_2,
    trigram_df,
    bigram_counts,
    vocabulary_size,
    preprocess_text
)
print(f"Trigram Probability for Sentence 2: {trigram_prob_2:.10e}")

```

--- Calculating Probabilities for Example Sentence 1: 'She was a woman of good sense.' ---

Unigram Probability for Sentence 1: 1.3231645770e-27

Bigram Probability for Sentence 1: 4.5756669350e-24

Trigram Probability for Sentence 1: 3.6222490553e-30

--- Calculating Probabilities for Example Sentence 2: 'The strange cat meowed loudly at the moon.' ---

Unigram Probability for Sentence 2: 4.1331803661e-40

Bigram Probability for Sentence 2: 2.0508221409e-33

Trigram Probability for Sentence 2: 1.6971095378e-35

Explanation of Sentence Probabilities

We calculated the probabilities for two example sentences using Unigram, Bigram (with add-one smoothing), and Trigram (with add-one smoothing) models.

Example Sentence 1: 'She was a woman of good sense.'

- **Unigram Probability:** 1.3231645770e-27
- **Bigram Probability:** 4.5756669350e-24
- **Trigram Probability:** 3.6222490553e-30

Observations for Sentence 1:

- The Bigram model yielded the highest probability for this sentence, followed by the Unigram model, and then the Trigram model. This suggests that the word sequences in this sentence are relatively common in bigram form but less so in trigram form, or that the independence assumption of the Unigram model is too loose for this specific sentence structure.
- The Unigram model, by considering words independently, assigns a probability based solely on individual word frequencies. The Bigram and Trigram models, incorporating context, tend to assign lower probabilities if specific word combinations are rare or unseen, even if individual words are common. The fact that the Trigram probability is lower than Bigram suggests that some trigrams in this sentence might be less frequent, despite its common English phrasing.

Example Sentence 2: 'The strange cat meowed loudly at the moon.'

- **Unigram Probability:** 4.1331803661e-40
- **Bigram Probability:** 2.0508221409e-33
- **Trigram Probability:** 1.6971095378e-35

Observations for Sentence 2:

- Similar to Sentence 1, the Bigram model provided a higher probability than the Trigram model, and both were significantly higher than the Unigram model. However, for this sentence, all probabilities are much lower than for Sentence 1, indicating that this sentence, especially words like 'strange' and 'meowed', or their combinations, are less frequent in the Jane Austen corpus.
- The significant difference in probabilities between Unigram and Bigram/Trigram models highlights the importance of context. The Unigram model struggles more with less common words because it doesn't account for their surrounding words, leading to a much lower overall probability. The higher-order models (Bigram and Trigram) can leverage the sequential information, potentially

assigning a higher probability if the word combinations are reasonable, even if individual words are rare. However, if combinations are rare or unseen, smoothing helps assign a non-zero, albeit very small, probability.

Explanation of Sentence Probabilities

We calculated the probabilities for two example sentences using Unigram, Bigram (with add-one smoothing), and Trigram (with add-one smoothing) models.

Example Sentence 1: 'She was a woman of good sense.'

- **Unigram Probability:** 1.3231645770e-27
- **Bigram Probability:** 4.5756669350e-24
- **Trigram Probability:** 3.6222490553e-30

Observations for Sentence 1:

- The Bigram model yielded the highest probability for this sentence, followed by the Unigram model, and then the Trigram model. This suggests that the word sequences in this sentence are relatively common in bigram form but less so in trigram form, or that the independence assumption of the Unigram model is too loose for this specific sentence structure.
- The Unigram model, by considering words independently, assigns a probability based solely on individual word frequencies. The Bigram and Trigram models, incorporating context, tend to assign lower probabilities if specific word combinations are rare or unseen, even if individual words are common. The fact that the Trigram probability is lower than Bigram suggests that some trigrams in this sentence might be less frequent, despite its common English phrasing.

Example Sentence 2: 'The strange cat meowed loudly at the moon.'

- **Unigram Probability:** 4.1331803661e-40
- **Bigram Probability:** 2.0508221409e-33
- **Trigram Probability:** 1.6971095378e-35

Observations for Sentence 2:

- Similar to Sentence 1, the Bigram model provided a higher probability than the Trigram model, and both were significantly higher than the Unigram model. However, for this sentence, all probabilities are much lower than for Sentence 1, indicating that this sentence, especially words like 'strange' and 'meowed', or their combinations, are less frequent in the Jane Austen corpus.
- The significant difference in probabilities between Unigram and Bigram/Trigram models highlights the importance of context. The Unigram model struggles more with less common words because it doesn't account for their surrounding words, leading to a much lower overall probability. The higher-order models (Bigram and Trigram) can leverage the sequential information, potentially assigning a higher probability if the word combinations are reasonable, even if individual words are rare. However, if combinations are rare or unseen, smoothing helps assign a non-zero, albeit very small, probability.

✓ Evaluate Models using Perplexity

Subtask:

Implement the perplexity metric and use it to evaluate the performance of the Unigram, Bigram, and Trigram language models on a test set (or a portion of the training set if a separate test set isn't available). Display the perplexity scores for each model.

Reasoning: I will define the `calculate_perplexity` function as instructed, which will compute the perplexity for Unigram, Bigram, and Trigram models given a list of evaluation tokens and model-specific parameters. This function will be crucial for evaluating the language models' performance.

```
import math

def calculate_perplexity(eval_tokens, model_name, unigram_df=None, total_words=None, vocabulary_size=None, bigram_df=None, word_log_prob_sum = 0.0
N = 0

if model_name == 'unigram':
    for token in eval_tokens:
        prob = get_unigram_probability(token, unigram_df, total_words, vocabulary_size)
        if prob > 0:
            log_prob_sum += math.log2(prob)
        N += 1
elif model_name == 'bigram':
    # Add a start token to eval_tokens for bigram context if not already present
```

```

# For perplexity, we evaluate P(w_i | w_{i-1}). The first word P(w_1) is often ignored or handled specially.
# If the tokens already include <s> and </s>, we can use them.
# Assuming eval_tokens already contains <s> and </s> for sentences within it.

# For each bigram, we need P(current_word | prev_word)
for i in range(1, len(eval_tokens)): # Start from 1 to form (w_{i-1}, w_i)
    prev_word = eval_tokens[i-1]
    current_word = eval_tokens[i]
    prob = get_bigram_probability(prev_word, current_word, bigram_df, word_counts, vocabulary_size)
    if prob > 0:
        log_prob_sum += math.log2(prob)
    N += 1

elif model_name == 'trigram':
    # Add '<s>' to eval_tokens to provide context for the first actual word in the eval_tokens
    # and for the second word as well.
    # If the tokens already include <s> and </s>, we use them.
    # For each trigram, we need P(current_word | prev_prev_word, prev_word)
    for i in range(2, len(eval_tokens)): # Start from 2 to form (w_{i-2}, w_{i-1}, w_i)
        prev_prev_word = eval_tokens[i-2]
        prev_word = eval_tokens[i-1]
        current_word = eval_tokens[i]
        prob = get_trigram_probability(prev_prev_word, prev_word, current_word, trigram_df, bigram_counts, vocabulary_size)
        if prob > 0:
            log_prob_sum += math.log2(prob)
        N += 1
    else:
        raise ValueError("Invalid model_name. Choose 'unigram', 'bigram', or 'trigram'.") 

if N == 0:
    return float('inf')

# Perplexity = 2 ^ (-1/N * sum(log2(P(wi))))
perplexity = 2 ** (-log_prob_sum / N)
return perplexity

print("Perplexity calculation function defined.")

Perplexity calculation function defined.

```

Reasoning: Now that the `calculate_perplexity` function is defined, I will split the `processed_tokens` into training and testing sets (80% for training, 20% for testing) to evaluate the models on unseen data, as specified in the instructions. This split will provide the `eval_tokens` for the perplexity calculations.

```

split_point = int(len(processed_tokens) * 0.8)

train_tokens = processed_tokens[:split_point]
test_tokens = processed_tokens[split_point:]

print(f"Total tokens: {len(processed_tokens)}")
print(f"Training tokens: {len(train_tokens)}")
print(f"Test tokens: {len(test_tokens)}")

Total tokens: 118711
Training tokens: 94968
Test tokens: 23743

```

Reasoning: I will calculate the perplexity for the Unigram, Bigram, and Trigram models using the `test_tokens` and the `calculate_perplexity` function, as instructed. Then I will print the perplexity score for each model to evaluate their performance.

```

print("\n--- Evaluating Perplexity on Test Set ---")

# Calculate Unigram Perplexity
unigram_perplexity = calculate_perplexity(
    test_tokens,
    model_name='unigram',
    unigram_df=unigram_df,
    total_words=total_words,
    vocabulary_size=vocabulary_size
)
print(f"Unigram Perplexity: {unigram_perplexity:.2f}")

# Calculate Bigram Perplexity

```

```

bigram_perplexity = calculate_perplexity(
    test_tokens,
    model_name='bigram',
    bigram_df=bigram_df,
    word_counts=word_counts,
    vocabulary_size=vocabulary_size
)
print(f"Bigram Perplexity: {bigram_perplexity:.2f}")

# Calculate Trigram Perplexity
trigram_perplexity = calculate_perplexity(
    test_tokens,
    model_name='trigram',
    trigram_df=trigram_df,
    bigram_counts=bigram_counts,
    vocabulary_size=vocabulary_size
)
print(f"Trigram Perplexity: {trigram_perplexity:.2f}")

```

```

--- Evaluating Perplexity on Test Set ---
Unigram Perplexity: 615.45
Bigram Perplexity: 1271.62
Trigram Perplexity: 3099.42

```

```
print(cleaned_text[:2000])
```

```
Is, he added to his wealth. To him therefore the succession to the Norland estate was not so really important as to his sisters;
```

Compare Model Performance

Subtask:

Analyze and compare the performance of the Unigram, Bigram, and Trigram models based on their perplexity scores. Discuss which N-gram model provides better predictions and why.

Compare Model Performance

Analysis of Model Performance based on Perplexity

We have calculated the perplexity scores for the Unigram, Bigram, and Trigram models on the test set:

- **Unigram Perplexity:** 615.45
- **Bigram Perplexity:** 1271.62
- **Trigram Perplexity:** 3099.42

Observations and Discussion:

Contrary to typical expectations where higher-order N-gram models usually perform better (i.e., have lower perplexity) due to capturing more context, our results show the Unigram model achieving the lowest perplexity, followed by Bigram, and then Trigram with the highest perplexity. A lower perplexity indicates a better language model, as it is more confident and surprised less by the test data.

This unexpected result suggests a few potential factors:

- Sparsity and Smoothing Impact:** As the order of the N-gram model increases, the number of unique N-grams (especially trigrams) grows significantly, leading to data sparsity issues. While add-one smoothing was applied, it might be overly aggressive for higher-order N-grams in this dataset. Add-one smoothing assigns a small, non-zero probability to all unseen N-grams. If many unique bigrams and trigrams in the test set were not present in the training data, applying smoothing to a large number of these unseen combinations could lead to a higher overall perplexity for the Bigram and Trigram models. This means the smoothed probabilities for observed higher-order N-grams might be significantly underestimated, or the baseline probabilities for unseen N-grams become too large when distributed across a very large vocabulary squared/cubed.
- Dataset Characteristics:** The specific nature of the Jane Austen text (formal, relatively constrained vocabulary, perhaps repetitive sentence structures) might mean that individual word frequencies (unigram) are highly indicative, and the added context from bigrams and trigrams introduces too much noise or too many rare sequences that are penalized heavily by smoothing.
- Preprocessing Impact:** The preprocessing step, especially the removal of all punctuation and numbers, might homogenize the text too much, leading to a loss of subtle contextual cues that higher-order models would typically benefit from. Additionally, the

inclusion of `<ss>` and `</ss>` tokens around every sentence could significantly inflate the vocabulary and the number of unique N-grams, exacerbating sparsity.

In this particular evaluation, the **Unigram model performed the best**, indicated by its significantly lower perplexity score. This suggests that for this specific dataset and the smoothing method used, the simple frequency of individual words was a more robust predictor than conditional probabilities derived from word sequences, likely due to severe data sparsity issues in the higher-order models despite smoothing.

Compare Model Performance

Analysis of Model Performance based on Perplexity

We have calculated the perplexity scores for the Unigram, Bigram, and Trigram models on the test set:

- **Unigram Perplexity:** 615.45
- **Bigram Perplexity:** 1271.62
- **Trigram Perplexity:** 3099.42

Observations and Discussion:

Contrary to typical expectations where higher-order N-gram models usually perform better (i.e., have lower perplexity) due to capturing more context, our results show the Unigram model achieving the lowest perplexity, followed by Bigram, and then Trigram with the highest perplexity. A lower perplexity indicates a better language model, as it is more confident and surprised less by the test data.

This unexpected result suggests a few potential factors:

1. **Sparsity and Smoothing Impact:** As the order of the N-gram model increases, the number of unique N-grams (especially trigrams) grows significantly, leading to data sparsity issues. While add-one smoothing was applied, it might be overly aggressive for higher-order N-grams in this dataset. Add-one smoothing assigns a small, non-zero probability to all unseen N-grams. If many unique bigrams and trigrams in the test set were not present in the training data, applying smoothing to a large number of these unseen combinations could lead to a higher overall perplexity for the Bigram and Trigram models. This means the smoothed probabilities for observed higher-order N-grams might be significantly underestimated, or the baseline probabilities for unseen N-grams become too large when distributed across a very large vocabulary squared/cubed.
2. **Dataset Characteristics:** The specific nature of the Jane Austen text (formal, relatively constrained vocabulary, perhaps repetitive sentence structures) might mean that individual word frequencies (unigram) are highly indicative, and the added context from bigrams and trigrams introduces too much noise or too many rare sequences that are penalized heavily by smoothing.
3. **Preprocessing Impact:** The preprocessing step, especially the removal of all punctuation and numbers, might homogenize the text too much, leading to a loss of subtle contextual cues that higher-order models would typically benefit from. Additionally, the inclusion of `<ss>` and `</ss>` tokens around every sentence could significantly inflate the vocabulary and the number of unique N-grams, exacerbating sparsity.

In this particular evaluation, the **Unigram model performed the best**, indicated by its significantly lower perplexity score. This suggests that for this specific dataset and the smoothing method used, the simple frequency of individual words was a more robust predictor than conditional probabilities derived from word sequences, likely due to severe data sparsity issues in the higher-order models despite smoothing.

Final Task

Subtask:

Provide a summary report of the entire lab, including the implementation details of each N-gram model, the evaluation results using perplexity, and the comparative analysis of model performance.

Summary:

Data Analysis Key Findings

- **Dataset Overview and Preprocessing:**
 - The analysis utilized Jane Austen's 'Sense and Sensibility' text, comprising approximately 668,777 characters.
 - Preprocessing involved converting text to lowercase, removing punctuation and numbers, tokenizing words, and adding "`(start)`" and "`(end)`" tokens to each sentence.

- The total number of processed tokens was 118,711, with a vocabulary size of 7,283 unique words.

- **N-gram Model Implementations:**

- **Unigram Model:** Calculated word frequencies (e.g., 'the' appeared 4,087 times) and their probabilities (e.g., 'the' had a probability of 0.0344).
- **Bigram Model:** Generated 118,710 bigrams, with 54,786 unique bigrams. Conditional probabilities were calculated using add-one (Laplace) smoothing, where $P(w_i|w_{i-1}) = (\text{Count}(w_{i-1}, w_i) + 1)/(\text{Count}(w_{i-1}) + \text{Vocabulary_Size})$. For instance, ('to', 'be') was the most frequent bigram.
- **Trigram Model:** Generated 118,709 trigrams, with 101,101 unique trigrams. Conditional probabilities were calculated using add-one smoothing, where $P(w_i|w_{i-2}, w_{i-1}) = (\text{Count}(w_{i-2}, w_{i-1}, w_i) + 1)/(\text{Count}(w_{i-2}, w_{i-1}) + \text{Vocabulary_Size})$. For example, ('i', 'am', 'sure') was the most frequent trigram with 68 occurrences and a smoothed probability of 0.009206.

- **Sentence Probability Calculations:** Functions were developed for each model to compute the probability of given sentences.

- For the sentence "She was a woman of good sense.":

- Unigram Probability: 1.323×10^{-27}
- Bigram Probability: 4.576×10^{-24}
- Trigram Probability: 3.622×10^{-30}

- For the sentence "The strange cat meowed loudly at the moon.":

- Unigram Probability: 4.133×10^{-40}
- Bigram Probability: 2.051×10^{-33}
- Trigram Probability: 1.697×10^{-35}

- **Perplexity Evaluation:**

- The data was split into an 80% training set (94,968 tokens) and a 20% test set (23,743 tokens).
- The perplexity scores on the test set were:

- Unigram Perplexity: 615.45
- Bigram Perplexity: 1271.62
- Trigram Perplexity: 3099.42

- **Comparative Analysis of Model Performance:**

- The Unigram model exhibited the best performance with the lowest perplexity (615.45), followed by the Bigram and Trigram models. This outcome is unexpected, as higher-order N-gram models typically perform better due to incorporating more context.
- This unexpected result is attributed to severe data sparsity issues in higher-order N-gram models, where the add-one smoothing technique, while ensuring non-zero probabilities, might be overly aggressive or insufficiently handle the vast number of unseen bigrams and trigrams in the test set.
- The specific characteristics of the Jane Austen corpus (formal, potentially repetitive language) and the comprehensive preprocessing (removing all punctuation and numbers, adding '⁇' tokens) may have exacerbated sparsity, making individual word frequencies (unigram) more robust predictors in this specific scenario.

Insights or Next Steps

- **Explore advanced smoothing techniques:** Investigate more sophisticated smoothing methods like Kneser-Ney or Witten-Bell interpolation, which are designed to handle data sparsity more effectively than simple add-one smoothing, especially for higher-order N-gram models.
- **Refine preprocessing strategies:** Experiment with alternative preprocessing steps, such as preserving some punctuation (e.g., sentence-ending marks) or using a different tokenization approach, to see if richer contextual cues can be maintained for higher-order models without excessively increasing sparsity.

