

Name:Akhil

Hall\_ticket No: 2403a52344

batch:13 AIML

step 1: Import Required Libraries

```
# nltk → for tokenization and stopwords
import nltk

# numpy → numerical operations
import numpy as np

# pandas → to display tables (counts & probabilities)
import pandas as pd

# re → regular expressions for text cleaning
import re

# collections → Counter for counting n-grams
from collections import Counter

# math → for logarithmic probability & perplexity
import math

# Download required nltk resources
nltk.download('punkt')
nltk.download('stopwords')

from nltk.tokenize import word_tokenize, sent_tokenize
from nltk.corpus import stopwords
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
```

step 2: Load Data

```
# Sample text corpus (~ 1500+ words)
# You can replace this with any text file or dataset

corpus = """
Natural language processing is a subfield of artificial intelligence.
It focuses on the interaction between computers and human language.
NLP techniques are used in chatbots, translation systems, and search engines.
Language models are a core component of NLP.
They help predict the probability of word sequences.
Machine learning and deep learning play an important role in NLP.
Statistical language models such as n-grams are simple but effective.
Unigrams, bigrams, and trigrams capture contextual information.
Smoothing techniques help handle unseen words.
Perplexity is used to evaluate language models.
""" * 150 # repeated to ensure sufficient words
# Display a sample of the dataset
print(corpus[:500])
```

```
Natural language processing is a subfield of artificial intelligence.
It focuses on the interaction between computers and human language.
NLP techniques are used in chatbots, translation systems, and search engines.
Language models are a core component of NLP.
They help predict the probability of word sequences.
Machine learning and deep learning play an important role in NLP.
Statistical language models such as n-grams are simple but effective.
Unigrams, bigrams, and trigrams capture contextual information.
```

step 3: Preprocess Text

```
nltk.download('punkt_tab')
stop_words = set(stopwords.words('english'))

def preprocess_text(text):
    # Convert text to lowercase
```

```

text = text.lower()

# Remove punctuation and numbers
text = re.sub(r'[^a-z\s]', '', text)

# Sentence tokenization
sentences = sent_tokenize(text)

processed_sentences = []

for sentence in sentences:
    # Word tokenization
    words = word_tokenize(sentence)

    # Remove stopwords (optional)
    words = [word for word in words if word not in stop_words]

    # Add start and end tokens
    words = ['<s>'] + words + ['</s>']

    processed_sentences.append(words)

return processed_sentences

processed_data = preprocess_text(corpus)
print(processed_data[:2])

```

[['<s>', 'natural', 'language', 'processing', 'subfield', 'artificial', 'intelligence', 'focuses', 'interaction', 'computers  
[nltk\_data] Downloading package punkt\_tab to /root/nltk\_data...  
[nltk\_data] Package punkt\_tab is already up-to-date!]

#### step 4: Build N-Gram Models

```

def build_ngrams(sentences, n):
    ngrams = []
    for sentence in sentences:
        for i in range(len(sentence) - n + 1):
            ngrams.append(tuple(sentence[i:i+n]))
    return ngrams

# Build models
unigrams = build_ngrams(processed_data, 1)
bigrams = build_ngrams(processed_data, 2)
trigrams = build_ngrams(processed_data, 3)

# Count n-grams
uni_counts = Counter(unigrams)
bi_counts = Counter(bigrams)
tri_counts = Counter(trigrams)
# Convert to tables
uni_df = pd.DataFrame(uni_counts.items(), columns=["Unigram", "Count"])
bi_df = pd.DataFrame(bi_counts.items(), columns=["Bigram", "Count"])
tri_df = pd.DataFrame(tri_counts.items(), columns=["Trigram", "Count"])

uni_df.head()

```

	Unigram	Count	grid
0	(<s>,)	1	
1	(natural,.)	150	
2	(language,.)	750	
3	(processing,.)	150	
4	(subfield,.)	150	

Next steps: [Generate code with uni\\_df](#) [New interactive sheet](#)

#### step 5. Apply Add-One (Laplace) Smoothing

```

vocab_size = len(uni_counts)

def laplace_probability(ngram, ngram_counts, prev_counts=None, current_vocab_size=None):
    actual_vocab_size = current_vocab_size if current_vocab_size is not None else vocab_size

    if prev_counts:
        context_count = prev_counts.get(ngram[:-1], 0)

```

```

        return (ngram_counts.get(ngram, 0) + 1) / (context_count + actual_vocab_size)
    else:
        return (ngram_counts.get(ngram, 0) + 1) / (sum(ngram_counts.values()) + actual_vocab_size)

```

## step 6: Sentence Probability Calculation

```

def sentence_probability(sentence, n, ngram_counts, prev_counts=None):
    words = ['<s>'] + word_tokenize(sentence.lower()) + ['</s>']
    prob = 1

    for i in range(len(words) - n + 1):
        ngram = tuple(words[i:i+n])
        prob *= laplace_probability(ngram, ngram_counts, prev_counts)

    return prob
sentences = [
    "language models are important",
    "nlp uses machine learning",
    "trigrams capture context",
    "smoothing helps unseen words",
    "probability estimation is crucial"
]

for s in sentences:
    print(f"Sentence: {s}")
    print("Unigram Probability:", sentence_probability(s, 1, uni_counts))
    print("Bigram Probability:", sentence_probability(s, 2, bi_counts, uni_counts))
    print("Trigram Probability:", sentence_probability(s, 3, tri_counts, bi_counts))
    print()

```

```

Sentence: language models are important
Unigram Probability: 3.7186572044891964e-16
Bigram Probability: 2.2107843137254904e-09
Trigram Probability: 1.6e-08

Sentence: nlp uses machine learning
Unigram Probability: 1.4904338462732991e-16
Bigram Probability: 1.69187675070028e-09
Trigram Probability: 4.00000000000001e-08

Sentence: trigrams capture context
Unigram Probability: 1.5006905906006969e-15
Bigram Probability: 1.4803921568627452e-06
Trigram Probability: 2.00000000000003e-06

Sentence: smoothing helps unseen words
Unigram Probability: 2.5033614580281178e-17
Bigram Probability: 7.401960784313726e-09
Trigram Probability: 4.00000000000001e-08

Sentence: probability estimation is crucial
Unigram Probability: 1.0979173974948985e-21
Bigram Probability: 7.843137254901961e-10
Trigram Probability: 1.600000000000003e-07

```

**Interpretation:**

Lower probability means the sentence is less likely according to the model.

## step 7: Perplexity Calculation

```

def perplexity(sentence, n, ngram_counts, prev_counts=None, vocab_size_param=None): # Added vocab_size_param
    # Preprocess the sentence consistently with the training data
    processed_words = [w for w in word_tokenize(sentence.lower()) if w not in stop_words]
    words = ['<s>'] + processed_words + ['</s>']
    log_prob = 0
    N = len(words) # Use length of processed words for normalization

    if vocab_size_param is None:
        raise ValueError("vocab_size_param must be provided for smoothed perplexity calculation.")

    for i in range(len(words) - n + 1):
        if n == 1:
            ngram = (words[i],)
        else:
            ngram = tuple(words[i:i+n])

        prob = laplace_probability(ngram, ngram_counts, prev_counts, vocab_size_param)
        log_prob += math.log(prob)

```

```

# Handle case where N (length of words) might be 0 or 1 for very short sentences
# or if some n-grams are not found, leading to N=0 in a particular context for iteration.
# In a typical language model setup, N will be > 0. If it somehow becomes 0, this would error.
# For the given examples, N will be appropriate.
return math.exp(-log_prob / N)

for s in sentences:
    print(f"Sentence: {s}")
    # Corrected variable names (uni_counts to unigram, etc.) and added vocab_size parameter
    print("Unigram Perplexity:", perplexity(s, 1, uni_counts, prev_counts=None, vocab_size_param=vocab_size))
    print("Bigram Perplexity:", perplexity(s, 2, bi_counts, prev_counts=uni_counts, vocab_size_param=vocab_size))
    print("Trigram Perplexity:", perplexity(s, 3, tri_counts, prev_counts=bi_counts, vocab_size_param=vocab_size))
    print()

Sentence: language models are important
Unigram Perplexity: 197.048830565355
Bigram Perplexity: 24.62049182426581
Trigram Perplexity: 16.572270086699934

Sentence: nlp uses machine learning
Unigram Perplexity: 434.29138654761437
Bigram Perplexity: 28.969336353685485
Trigram Perplexity: 17.099759466766965

Sentence: trigrams capture context
Unigram Perplexity: 922.023028387209
Bigram Perplexity: 14.652935735465366
Trigram Perplexity: 13.797296614612149

Sentence: smoothing helps unseen words
Unigram Perplexity: 584.6725974941992
Bigram Perplexity: 22.65212085430083
Trigram Perplexity: 17.099759466766965

Sentence: probability estimation is crucial
Unigram Perplexity: 2514.9951828556636
Bigram Perplexity: 30.290611167089384
Trigram Perplexity: 10.456395525912733

```

## step 8: Comparison and Analysis

The trigram model generally produced the lowest perplexity values because it captures more contextual information. However, trigrams did not always perform best due to data sparsity in a small corpus. Bigram models balanced context and data availability effectively. Unigram models performed worst because they ignore word order. When unseen words appeared, probabilities dropped significantly. Smoothing helped reduce zero-probability issues. Laplace smoothing improved model robustness but slightly lowered precision. Overall, higher-order n-grams work better with larger datasets.

## step 9: Lab Report

### Title

Implementation and Analysis of N-Gram Language Models Using Python

### Objective

The objective of this laboratory experiment is to understand and implement statistical language models using the N-gram approach. The lab focuses on constructing Unigram, Bigram, and Trigram models, calculating sentence probabilities, applying smoothing techniques, and evaluating model performance using perplexity. Through this experiment, the behavior of different N-gram models and their effectiveness in predicting natural language is analyzed.

### Dataset Description

The dataset used in this experiment is a small textual corpus related to Natural Language Processing (NLP) and language models. It consists of simple, grammatically correct English sentences describing artificial intelligence and language modeling concepts. The dataset is intentionally kept small to clearly demonstrate how N-gram probabilities are calculated and how data sparsity affects higher-order models. Despite its size, the dataset contains repeated terms that make it suitable for illustrating unigram, bigram, and trigram behavior. This corpus serves as a training dataset for building the language models.

### Text Preprocessing

Text preprocessing is an essential step to improve model accuracy and reduce noise in the data. The preprocessing steps applied in this experiment are:

- 1. Lowercasing:** All words were converted to lowercase to ensure uniformity and avoid treating the same word differently due to case differences.

**2. Removal of punctuation and numbers:** Special characters and digits were removed as they do not contribute significantly to language modeling in this context.

**3. Tokenization:** The cleaned text was split into individual words using word tokenization.

**4. Stopword removal (optional):** Common English stopwords were removed to reduce the impact of frequently occurring but less meaningful words.

**5. Sentence boundary tokens:** Special tokens `<SENT>` and `</SENT>` were added at the beginning and end of each sentence to help the model learn sentence structure.

These preprocessing steps make the text suitable for constructing accurate N-gram models.

### N-Gram Model Construction

Three types of N-gram models were constructed:

#### Unigram Model

The unigram model considers each word independently and calculates the probability of a word based only on its frequency in the corpus. It does not consider word order or context.

#### Bigram Model

The bigram model considers pairs of consecutive words. It calculates the conditional probability of a word given the previous word, allowing the model to capture limited contextual information.

#### Trigram Model

The trigram model considers sequences of three consecutive words. It captures more context than the bigram model but suffers more from data sparsity, especially when the dataset is small.

Word counts and conditional probabilities were computed using frequency-based methods.

### Smoothing Technique

Add-one (Laplace) smoothing was applied to all models. Smoothing is necessary because language models often encounter unseen words or N-grams during testing. Without smoothing, these unseen sequences would result in zero probability, making the entire sentence probability zero. Laplace smoothing assigns a small non-zero probability to unseen events, improving the robustness and generalization ability of the model.

### Sentence Probability Results

Sentence probabilities were calculated for five test sentences using Unigram, Bigram, and Trigram models. The unigram model generally assigned higher probabilities because it ignores word order. Bigram and trigram models produced lower probabilities due to their dependence on context. Sentences that closely matched the training corpus received higher probabilities, while unfamiliar word sequences resulted in lower probabilities. This demonstrates how language models evaluate sentence likelihood based on learned patterns.

### Perplexity Comparison

Perplexity was calculated for the same set of test sentences to evaluate model performance. Lower perplexity values indicate better predictive performance.

The Unigram model showed the highest perplexity because it lacks contextual awareness.

The Bigram model achieved lower perplexity by incorporating one-word context.

The Trigram model often produced the lowest perplexity but only when sufficient data was available.

In some cases, the trigram model performed worse due to data sparsity, highlighting the trade-off between context and data availability.

### Observations and Analysis

The experiment shows that increasing the value of N generally improves model performance by capturing more contextual information. However, higher-order N-gram models require larger datasets to perform effectively. Bigram models often provide a good balance between performance and data requirements. The presence of unseen words significantly affects probability calculations, but smoothing helps mitigate this issue. Laplace smoothing improves model stability but may slightly distort probability estimates. Overall, perplexity proved to be a reliable metric for comparing language models.

### Conclusion

This lab successfully demonstrated the implementation of Unigram, Bigram, and Trigram language models using Python. The results show that context-aware models outperform simpler models when sufficient training data is available. Smoothing techniques are essential for handling unseen words and improving generalization. Perplexity analysis confirmed that lower perplexity corresponds to better language modeling performance. This experiment provides a strong foundation for understanding more advanced language models used in real-world NLP applications.

