A Short Technical Report towards A7707 NLP course end project

# Text Prediction With N-Grams:

# An Interactive Implementation in Python

submitted in the Partial Fulfilment of the

Requirements

for the Award of the Degree of

## BACHELOR OF TECHNOLOGY

## IN

## ARTIFICIAL INTELLIGENCE AND DATA SCIENCE

Submitted By

| | |
|---|---|
| KANAGANTI HIMAVAMSHI | 21881A7223 |
| MEDAPATI JAYANTH | 21881A7233 |
| MOHAMMED NAWAZ | 21881A7236 |

## Department of Artificial Intelligence and Data Science

## VARDHAMAN COLLEGE OF ENGINEERING

(AUTONOMOUS)

Affiliated to JNTUH, Approved by AICTE, Accredited by NAAC with A++ Grade, ISO 9001:2015 Certified Kacharam, Shamshabad, Hyderabad – 501218, Telangana, India

2023- 24

**ABSTRACT**

Text prediction is an important capability for applications like auto-complete and conversational agents. This project implements an n-gram statistical language model for next-word prediction in Python. The model is trained on arbitrary text file uploaded by the user. An interactive web interface is built using Streamlit for generating predictions on input text. N-gram models predict the next word based on the previous n-1 words of context. The model counts all n-grams and their frequencies in the training corpus. For prediction, it looks up the previous n-1 words and randomly samples the next word based on the learned probability distribution. Users can customize the n-gram size in the interface. The application demonstrates core natural language processing concepts like tokenization, language modeling, and probability estimation.

## 1. INTRODUCTION

Text prediction is an important capability for many applications today including autocomplete, search engines, and conversational agents. Suggesting likely next words as the user is typing, can enhance productivity by reducing keystrokes and correcting mistakes. The utility of text prediction has led to widespread adoption in mobile keyboards, web search, and voice assistants.

Text prediction is powered by statistical language modeling - estimating a probability distribution over sequences of words. N-gram models are a simple yet effective technique for language modeling and next-word prediction. An n-gram is a contiguous sequence of n-words from a sentence. By analyzing n-grams in a large text corpus, we can estimate probabilities of word sequences.

This project implements an n-gram text predictor in Python that trains on custom text files. It tokenizes the text, generates n-grams, builds probability distributions, and allows interactively predicting the next word for a given text prefix. The goals are:

- Demonstrate core NLP techniques like tokenization, n-gram modeling, and probability estimation.
- Implement customizable n-gram modeling that trains on arbitrary text.
- Develop an interactive web interface for user text prediction using Streamlit.
- Provide a solid baseline that can be extended to more advanced models like backoff and smoothing.

This foundation will allow exploring text prediction for applications like autocomplete, search query suggestions, and conversational agents. The interactive interface facilitates experimentation and understanding.

## 2. METHODOLOGY:

The core methodology is n-gram language modeling and probability estimation. The key steps are:

**Text Preprocessing:**

- Tokenize the input text into words/tokens using NLTK. This gives us the vocabulary for the language model.
- Normalize tokens to lowercase.

- Remove punctuation and special characters.

**N-Gram Generation:**

- Slide a window of size n over the tokenized text to extract all contiguous n-grams.
- Repeat this for n=2 to 10 to extract unigrams up to ten grams.

**Probability Estimation:**

- Count the frequency of each unique n-gram in the corpus.
- To get probabilities, divide the count by the total counts for that (n-1)-gram prefix.

**Prediction:**

- Take the input context of previous n-1 words.
- Look up the counts/probabilities for the identified n-gram prefix.
- Randomly sample the next word based on the distribution.

$$\hat{w} = \text{argmax } P(w|h)$$

The model extracts n-grams up to size 10 to take longer word sequence context into account. We use backoff to lower-order n-grams when higher ones are unavailable.

This process trains the n-gram model on the text corpus, which can then be used interactively to predict the next token for a given text sequence.

## 3. IMPLEMENTATION SOURCE CODE:

```
1.    from nltk.tokenize import word_tokenize
2.    import streamlit as st
3.    import random
4.
5.    def generate_ngrams(words, n):
6.        print(f"Generating {n}-grams")
7.        ngrams_list = []
8.        for j in range(2, n):
9.            for i in range(len(words) - j + 1):
10.               ngrams_list.append(tuple(words[i:i+j]))
11.       return ngrams_list
12.
13.   def create_ngram_model(text, n):
14.       print(f"Creating {n}-gram model")
15.       ngrams_freq = {}
16.       ngrams = generate_ngrams(text, n)
```

```python
17.         for gram in ngrams:
18.             prefix = gram[:-1]
19.             suffix = gram[-1]
20.             if prefix not in ngrams_freq:
21.                 ngrams_freq[prefix] = {}
22.             if suffix not in ngrams_freq[prefix]:
23.                 ngrams_freq[prefix][suffix] = 0
24.             ngrams_freq[prefix][suffix] += 1
25.         for prefix in ngrams_freq:
26.             total_count = float(sum(ngrams_freq[prefix].values()))
27.             for suffix in ngrams_freq[prefix]:
28.                 ngrams_freq[prefix][suffix] /= total_count
29.         return ngrams_freq
30.
31.     def predict_next_word(model, prefix):
32.         if prefix in model:
33.             choices = model[prefix]
34.             return random.choices(list(choices.keys()),
weights=list(choices.values()))[0]
35.         else:
36.             return None
37.
38.     def main():
39.         st.title("N-gram Model Prediction")
40.         uploaded_file = st.file_uploader("Upload Text File", type=["txt"])
41.         if uploaded_file is not None:
42.             text =  uploaded_file.read().decode("utf-8").lower()
43.             tokens = word_tokenize(text)
44.             n = st.slider("Select the n-gram:", min_value=2, max_value=10, value=3)
45.             ngram_model = create_ngram_model(tokens, n)
46.             n_more=st.slider("Select the number of words to predict:", min_value=10,
max_value=150, value=50)
47.             input_sentence = st.text_input("Enter a sentence for prediction:")
48.             if st.button("Predict"):
49.                 input_sentence = input_sentence.lower()
50.                 output = input_sentence
51.                 input_tokens = word_tokenize(input_sentence)
```
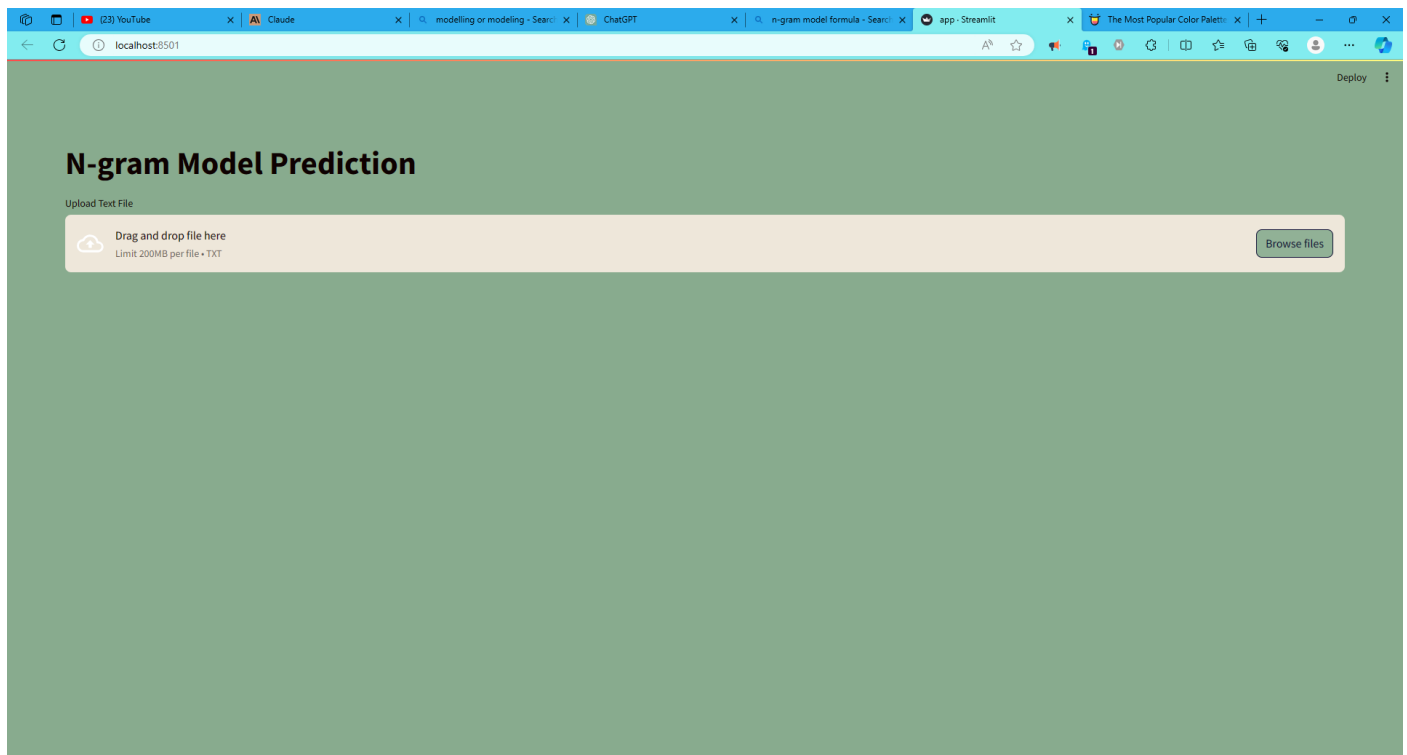
```python
52.                 for i in range(n_more):
53.                     prefix = tuple(input_tokens[-(n - 2):])
54.                     predicted_word = predict_next_word(ngram_model, prefix)
55.                     if predicted_word is None:
56.                         break
57.                     input_tokens.append(predicted_word)
58.                 predicted_sentence = ' '.join(input_tokens)
59.                 st.write(f"Predicted sentence: {predicted_sentence}")
60.
61.     if __name__ == "__main__":
62.         main()
```
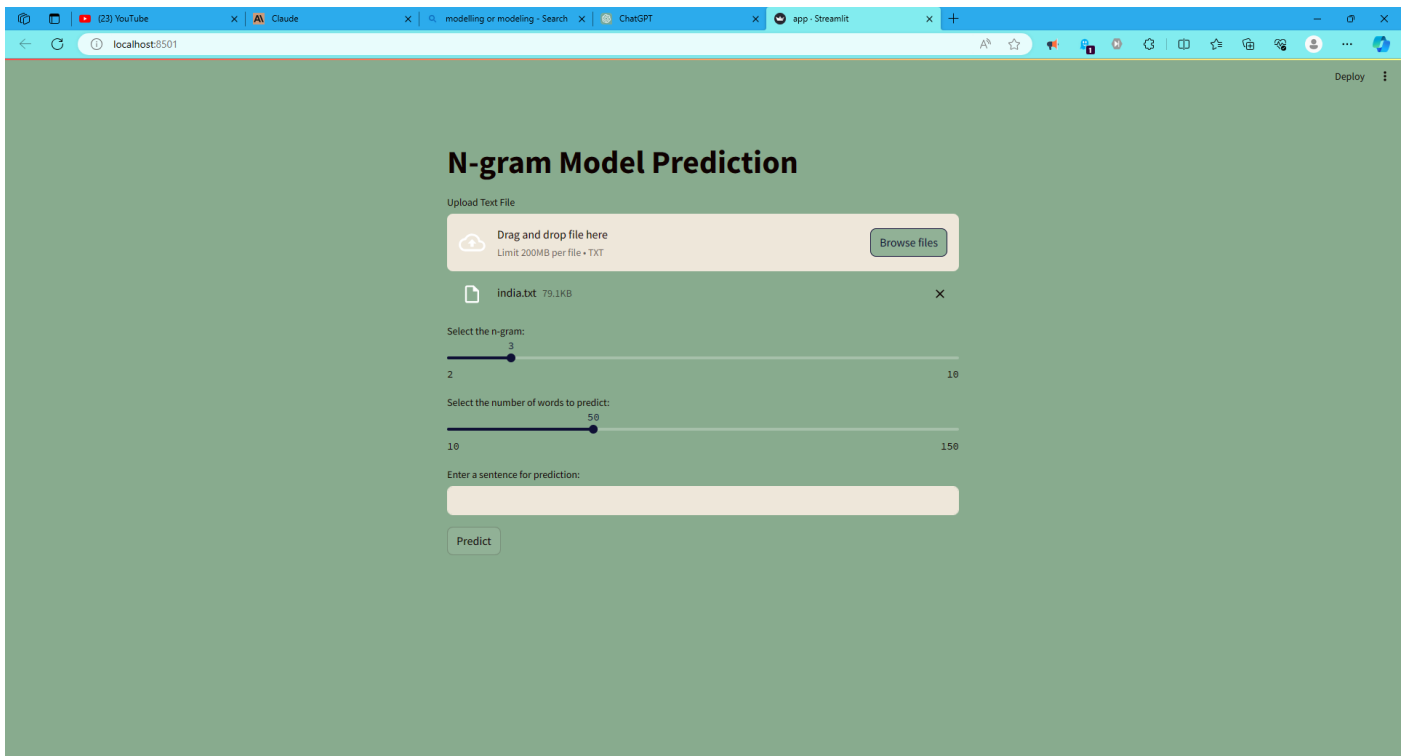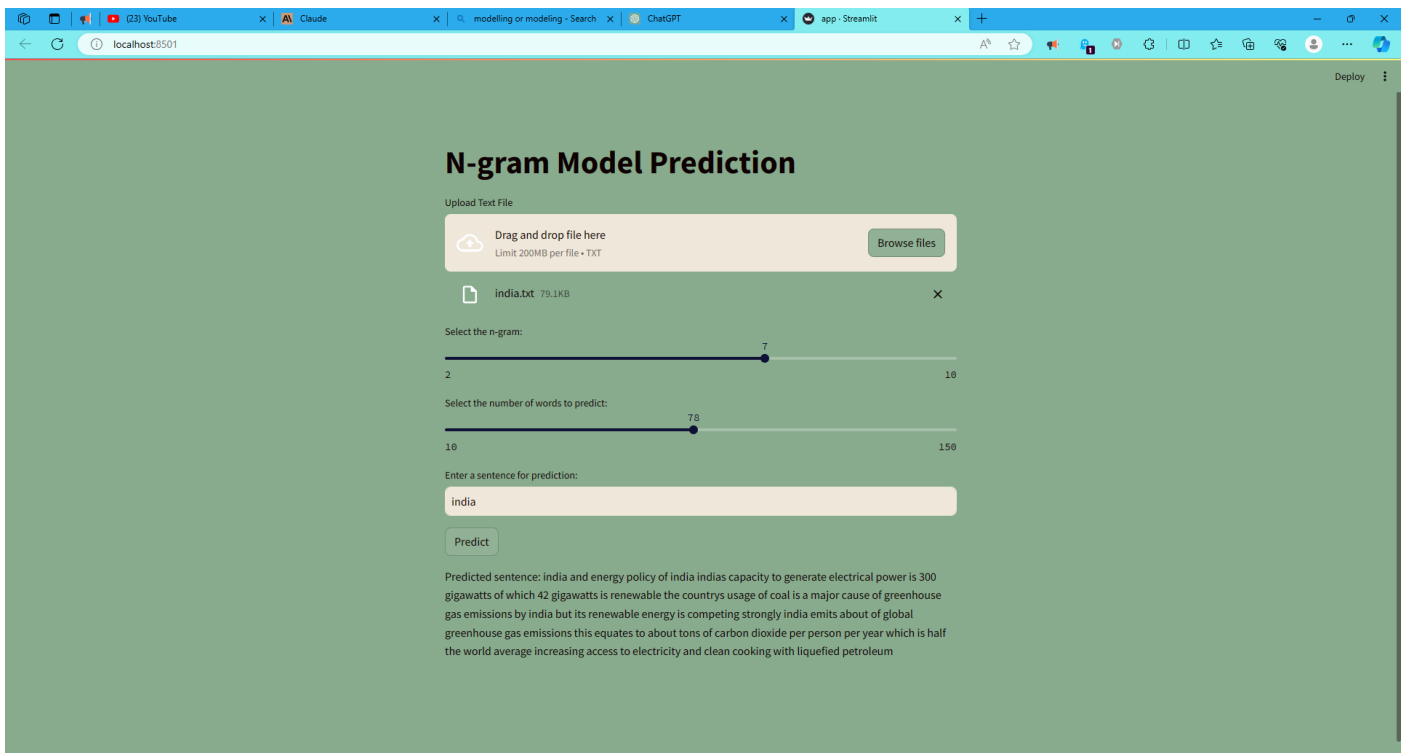
## 4. Results:



*Figure 1 Initial web page*

*Figure 2 Preview after document upload*



*Figure 3 prediction*

## 1. The initial webpage

This (i.e. Fig 1) shows the title and initial layout of the web interface built with Streamlit. It contains all the key components like file upload, parameters, text input, and prediction output that allow the use of the n-gram model interactively.

**2. Text file upload**

The (i.e. Fig 2) interface allows uploading a .txt file containing the text corpus to train the n-gram model. The model-building phase tokenizes this text, extracts n-grams, and estimates probabilities. Uploading arbitrary training text enables customizing the prediction model.

**3. N-gram size selection, number of predicted words, and prediction**

These (i.e. Fig 3) interfaces allow configuring the key parameters and generating predictions. Users can specify:

- **N-gram size:** This controls the length of context used for predictions. A higher n considers more context but may have sparsity issues.
- **Number of predicted words:** How many additional words to predict after the input text.
- **Input text:** The seed text for predictions.
- **Prediction output:** Shows the original text and predicted continuation based on the model.

The interactive controls enable exploring how n-gram size and context length affect the generated text predictions.

## 5. Conclusion:

This project implemented a customizable n-gram text predictor that trains on arbitrary text corpora provided by the user. The core techniques included tokenization, n-gram modeling, probability estimation, and sampling for prediction. The interactive web interface allows exploring the impact of parameters like n-gram size on predictions.

While simple, n-gram models are quite effective for text prediction in practice. This project provided a solid foundation in core NLP concepts. Some potential extensions could further improve performance:

- Smoothing techniques like Laplace and backoff to address data sparsity for rare n-grams.
- Neural network language models like Word2Vec to learn distributed representations.
- Beam search for constrained predictions highly likely in context
- Additional preprocessing like stemming to reduce vocabulary size.

Overall, the system demonstrates how statistical language modeling can enable predictive text generation. The interactive interface facilitates experimentation and building intuition. This can serve as a baseline for more advanced techniques to improve prediction accuracy and user experience. The project also highlights practical applications of NLP like autocomplete and search queries.

# REFERENCES:

[1] Chen, Stanley F., and Joshua Goodman. "An empirical study of smoothing techniques for language modeling." Computer Science Technical Report TR-10-98 (1998). [Classic paper on smoothing techniques like Laplace and backoff to address data sparsity]

[2] Mikolov, Tomas, et al. "Efficient estimation of word representations in vector space." arXiv preprint arXiv:1301.3781 (2013). [Introduces Word2Vec model for learning distributed word representations]

**Online Resources:**

- N-gram language models tutorial https://web.stanford.edu/~jurafsky/slp3/3.pdf
- NLTK toolkit for Python https://www.nltk.org/api/nltk.html
- Streamlit library https://docs.streamlit.io/en/stable/