

```
1 from preprocessing import Preprocessing
2 from ngrams import NGrams
3 from probability import Probability
4 from random import shuffle
5
6
7 def main():
8     # Set the dataset path
9     dataset_path = "../data/berkeley_restaurant_dataset.txt"
10
11     # Create an instance of the Preprocessing class with the dataset path
12     data = Preprocessing(dataset_path)
13
14     # Preprocess the dataset
15     data.preprocess()
16
17     # Grab the preprocessed sentences to print them
18     preprocessed_data = data.filtered_sentences
19
20     # Create an instance of the NGrams class with the preprocessed data
21     ngrams = NGrams(preprocessed_data)
22
23     # Generate unigrams, bigrams, and trigrams
24     ngrams.generate_ngrams()
25
26     # Access the generated ngrams
27     unigrams = ngrams.unigrams
28     bigrams = ngrams.bigrams
29     trigrams = ngrams.trigrams
30
31     # For testing - Print some example ngrams
32     # print("Unigrams:", unigrams[:10])
33     # print("Bigrams:", bigrams[:10])
34     # print("Trigrams:", trigrams[:10])
35
36     # Print the requested statistics on word count, vocabulary size,
```

```

37     # and total number of sentences
38     ngrams.print_statistics()
39
40     print("""
41
42         PART 3
43         -----
44
45         Read the chapter on N-grams and generate figures 4.1 and 4.2
46         for bigram
47         counts. The figures do not have to be exact.
48
49     """)
50     print("----- Bigram Count Table - Figure 4.1 -----\\n")
51     ngrams.bigram_count_table()
52
53     print("\\n----- Bigram Probability Table - Figure 4.2 -----\\n")
54     ngrams.bigram_probability_table()
55
56     print("""
57
58         PART 4 & 5
59         -----
60
61         Calculate the joint probability for at least five sentences (with vocabulary in the dataset)
62         using bigrams.
63         Repeat step 2 using trigrams. Observe if the estimates have changed.
64
65     """)
66     print("\\n----- Calculate bigram and trigram probabilities -----\\n")
67
68     probability = Probability(ngrams.bigrams, ngrams.trigrams)
69     bigram_probs = probability.bigram_probabilities()
70     trigram_probs = probability.trigram_probabilities()
71
72     sentences = data.filtered_sentences

```

```
72
73     shuffle(sentences)
74     examples = []
75     for sentence in sentences:
76         if len(sentence[1:]) > 2:
77             examples.append(sentence)
78
79     for sentence in examples[:5]:
80         if len(sentence[1:]) > 2:
81             print(f"Sentence: {' '.join(sentence[1:])}")
82             print(f"Bigram Probability: {bigram_probs[(sentence[1], sentence[2])]}")
83             print(f"Trigram Probability: {trigram_probs[(sentence[1], sentence[2], sentence[3])]}")
84             print()
85
86
87 if __name__ == "__main__":
88     main()
89
90
```

```
1 import random
2 import nltk
3 from collections import Counter
4 from tabulate import tabulate
5
6
7 class NGrams:
8     def __init__(self, preprocessed_data):
9         self.preprocessed_data = preprocessed_data
10        self.unigrams = []
11        self.bigrams = []
12        self.trigrams = []
13
14    def generate_unigrams(self):
15        for sentence in self.preprocessed_data:
16            for token in sentence[1:]:
17                self.unigrams.append(token)
18
19    def generate_bigrams(self):
20        for sentence in self.preprocessed_data:
21            bigrams = list(nltk.bigrams(sentence[1:]))
22            self.bigrams.extend(bigrams)
23
24    def generate_trigrams(self):
25        for sentence in self.preprocessed_data:
26            trigrams = list(nltk.trigrams(sentence[1:]))
27            self.trigrams.extend(trigrams)
28
29    def generate_ngrams(self):
30        self.generate_unigrams()
31        self.generate_bigrams()
32        self.generate_trigrams()
33
34    def count_words(self):
35        return len(self.unigrams)
36
```

```

37     def count_vocabulary(self):
38         return len(set(self.unigrams))
39
40     def count_sentences(self):
41         return len(self.preprocessed_data)
42
43     def print_statistics(self):
44         word_count = self.count_words()
45         vocab_size = self.count_vocabulary()
46         sentence_count = self.count_sentences()
47
48         print("""
49         PART 2
50         -----
51
52         - Count the words
53         - Report the size of the vocabulary
54         - report the number of sentences in the dataset
55
56         """)
57
58         print(f"WORD COUNT - Total number of words in the dataset:"
59               f" {word_count}")
60         print(f"VOCAB SIZE - Total number of unique words in the dataset:"
61               f" {vocab_size}")
62         print(f"SENTENCE COUNT - Number of sentences in the dataset:"
63               f" {sentence_count}")
64
65     def bigram_count_table(self, num_words = 8):
66         # Get a random set of unique words < -- Original goal
67         # commented out to use the words from the mentioned text
68         # unique_words = list(set(self.unigrams))
69         # selected_words = random.sample(unique_words, num_words)
70         selected_words = ['i', 'want', 'to', 'eat', 'chinese', 'food',
71                           'lunch', 'spend']
72

```

```
73     # Calculate bigram counts
74     bigram_counts = Counter(self.bigrams)
75
76     # Create a matrix with bigram counts
77     matrix = []
78     for w1 in selected_words:
79         row = []
80         for w2 in selected_words:
81             row.append(bigram_counts[(w1, w2)])
82         matrix.append(row)
83
84     # Display the matrix in a tabular format
85     headers = [''] + selected_words
86     table = tabulate(
87         matrix, headers = headers, showindex = selected_words,
88         tablefmt = "grid", numalign = "right"
89     )
90     print(table)
91
92     def bigram_probability_table(self, num_words = 8):
93         # Get a random set of unique words <-- Original goal
94         # commented out to use the words from the mentioned text
95         # unique_words = list(set(self.unigrams))
96         # selected_words = random.sample(unique_words, num_words)
97         selected_words = ['i', 'want', 'to', 'eat', 'chinese', 'food',
98                          'lunch', 'spend']
99
100        # Calculate unigram counts
101        unigram_counts = Counter(self.unigrams)
102
103        # Calculate bigram counts
104        bigram_counts = Counter(self.bigrams)
105
106        # Create a matrix with bigram probabilities
107        matrix = []
108        for w1 in selected_words:
```

```
109         row = []
110         for w2 in selected_words:
111             bigram_count = bigram_counts[(w1, w2)]
112             unigram_count = unigram_counts[w1]
113             probability = bigram_count / unigram_count if unigram_count > 0 else 0
114             row.append(probability)
115         matrix.append(row)
116
117     # Display the matrix in a tabular format
118     headers = [''] + selected_words
119     table = tabulate(
120         matrix, headers = headers, showindex = selected_words,
121         tablefmt = "grid", floatfmt = ".4f", numalign = "right"
122     )
123     print(table)
124
```

```
1 import nltk
2 from collections import defaultdict, Counter
3
4
5 class Probability:
6     def __init__(self, bigrams, trigrams):
7         self.bigrams = bigrams
8         self.trigrams = trigrams
9
10    def count_bigrams(self):
11        bigram_counts = defaultdict(Counter)
12        for bigram in self.bigrams:
13            bigram_counts[bigram[0]][bigram[1]] += 1
14
15        return bigram_counts
16
17    def bigram_probabilities(self):
18        bigram_counts = self.count_bigrams()
19
20        bigram_probs = {}
21        for first_word in bigram_counts:
22            total_count = sum(bigram_counts[first_word].values())
23            for second_word in bigram_counts[first_word]:
24                bigram_probs[(first_word, second_word)] = bigram_counts[first_word][second_word] /
total_count
25
26        return bigram_probs
27
28    def count_trigrams(self):
29        trigram_counts = defaultdict(Counter)
30        for trigram in self.trigrams:
31            trigram_counts[(trigram[0], trigram[1])][trigram[2]] += 1
32
33        return trigram_counts
34
35    def trigram_probabilities(self):
```



```
36         bigram_counts = self.count_bigrams()
37         trigram_counts = self.count_trigrams()
38
39         trigram_probs = {}
40         for bigram in trigram_counts:
41             for third_word in trigram_counts[bigram]:
42                 trigram_probs[bigram + (third_word,)] = trigram_counts[bigram][third_word] /
bigram_counts[bigram[0]][bigram[1]]
43
44         return trigram_probs
```

```

1 import nltk
2 from nltk.tokenize import word_tokenize, sent_tokenize
3
4
5 class Preprocessing:
6     def __init__(self, dataset_path, disfluencies=None):
7         self.dataset_path = dataset_path
8         self.raw_text = ""
9         self.sentences = []
10        self.filtered_sentences = []
11        if disfluencies is None:
12            self.disfluencies = ["uh", "uhm"]
13        else:
14            self.disfluencies = disfluencies
15
16    def read_data(self):
17        with open(self.dataset_path, "r") as file:
18            self.raw_text = file.read()
19
20    def remove_punctuation(self):
21        self.raw_text = self.raw_text.replace(" ", "")
22
23    def tokenize_sentences(self):
24        self.sentences = sent_tokenize(self.raw_text)
25        # Print Raw Data
26        print("""
27        PART 1
28        -----
29
30        Load and preprocess the dataset provided:
31        - Tokenize the text, keeping only actual words while removing disfluencies such as "uh"
32        and "uhm"
33        - Add special tokens to indicate the beginning of each sentence
34
35        """)
36        print("\n----- Sampled Raw Data - START-----\n")

```

```
37     for i in range(5):
38         print(self.sentences[i])
39     print("\n----- END ----- \n")
40
41     def tokenize_words(self):
42         self.filtered_sentences = [word_tokenize(sentence) for sentence in self.sentences]
43
44     def contains_non_alpha_chars(self, token):
45         for char in token:
46             if not char.isalpha():
47                 return True
48         return False
49
50     def remove_disfluencies(self):
51         for i, sentence in enumerate(self.filtered_sentences):
52             self.filtered_sentences[i] = [
53                 token.strip("_") for token in sentence
54                 if token.lower() not in self.disfluencies
55                 and not self.contains_non_alpha_chars(token)
56             ]
57
58     def add_start_tokens(self, start_token):
59         for sentence in self.filtered_sentences:
60             sentence.insert(0, start_token)
61
62     print("\n----- Sampled Processed Data - START----- \n")
63     for i in range(5):
64         print(self.filtered_sentences[i])
65     print("\n----- END ----- \n")
66
67     def preprocess(self, start_token="</s>"):
68         nltk.download("punkt")
69         self.read_data()
70         self.remove_punctuation()
71         self.tokenize_sentences()
72         self.tokenize_words()
```

```
73         self.remove_disfluencies()
74         self.add_start_tokens(start_token)
75
```