

February 14, 2024



Department of Computer Science and Engineering
Islamic University of Technology (IUT)
A subsidiary organ of OIC

Natural Language Processing, Assignment 1

Azmayen Fayek Sabil, 190042122

Contents

1	Task 1	3
1.1	Regular expressions	3
1.2	Basic Text Cleansing	3
1.3	Sentence Extractor	3
1.3.1	Explanation of the Regular Expression	4
1.4	Sentence Tokenizer	4
1.4.1	Regular Expression for Sentence Tokenization	5
1.4.2	Tokenizing the Entire Text	5
1.4.3	Tokenizing Matched Sentences	5
1.5	Compare Works	6
1.6	Word Frequency Distribution	6
1.7	Zipf's Law Validation	8
1.7.1	Analysis of Zipf's Law Conformance	8
1.7.2	Observations	8
1.7.3	Takeaway on Zipf's law	10
1.8	Comments on Statistical Similarities or Differences	10
1.9	Top 10 Words and Comparisons	10
1.10	Conclusion	10
2	Task 2	11
2.1	Text normalization	11
2.2	Tokenization with NLTK	11
2.3	Byte Pair Encoding	11
2.3.1	Introduction	11
2.3.2	Implementation	12
2.3.3	Explanation	13
2.3.4	Evaluation of Byte Pair Encoding	13
2.3.5	Plots	14
2.3.6	Observations	15
2.3.7	Results	15
3	Task 3	16
3.1	Edit Distance Measures	16
3.1.1	Hamming Distance	16
3.1.2	Levenshtein Distance	16
3.1.3	Needleman-Wunsch Algorithm	16
3.1.4	Smith-Waterman Algorithm	16
3.1.5	Code	16
3.2	Edit Distance Measures Output	17
3.2.1	Hamming Distance	17
3.2.2	Levenshtein Distance	17
3.2.3	Needleman-Wunsch Distance	18
3.2.4	Smith-Waterman Distance	18
3.3	Token-Level Distance Measures	18
3.3.1	Jaccard Index	18
3.3.2	Code Execution	18
3.3.3	Observations	19
3.4	Towards a More Nuanced Edit Distance Cost Metric	19
3.4.1	Implementation Explanation	19

4	Task 4	21
4.1	N-gram language models	21
4.1.1	Skip-Grams	21
4.2	Language Models	22
4.2.1	Building a Language Model	22
4.2.2	Converting Counts to Probabilities	22
4.2.3	Sampling from the Language Model	23
4.2.4	Generating New Sentences	23
4.3	Generating Philosophical Conversation	23

1 Task 1

1.1 Regular expressions

Regular expressions is basically used in text processing. We can use regular expression to find, extract patterns from a pile of text. Mainly it is used for searching, matching, and manipulating textual data. Here are some common regular expressions used in NLP:

1.2 Basic Text Cleansing

When working with downloaded books, it's common to perform basic text cleansing to ensure consistency and ease of analysis. Here's an example using Python's 're' module:

```
1 import re
2
3 def basic_text_cleansing(text):
4     # Replacing various forms of whitespace with a single space
5     text = re.sub('\s+', ' ', text)
6
7     # Converting the text to lowercase
8     text = text.lower()
9
10    # Replacing common abbreviations with their full form
11    text = re.sub("it's", 'it is', text)
12    text = re.sub("don't", 'do not', text)
13    text = re.sub("Don't", 'Do not', text)
14    text = re.sub("I'm", 'I am', text)
15    # text = re.sub("you are", 'you\'re', text)
16
17    return text
18
19 cleansed_mySite_text = basic_text_cleansing(mySite.text)
```

In this example, the *basic_text_cleansing(text)* function takes the text of the downloaded book and performs the following operations using 're.sub':

1. Replaces various forms of whitespace with a single space using the regular expression `\s+`.
2. Converts all text to lowercase using 'lower()' method.
3. Replaces the abbreviation like "it's" with its full form "it is".

We can also extend the function by adding more 're.sub' calls for other common abbreviations or patterns that need to be replaced in your specific text data.

1.3 Sentence Extractor

We can use regular expression to extract sentences that matches a certain expression.

This Python code below demonstrates the implementation of a regular expression to extract complete sentences containing matches for the term `(P|p)hilo(?!nous)[a-z]+`. The code processes three philosophical works: "Beyond Good and Evil," "The Problems of Philosophy," and a self-chosen philosophy work which is "Man and Superman: A Comedy and a Philosophy".

```
1 import requests
2 import re
3
4 # URLs for the philosophical works
5 website1 = 'https://www.gutenberg.org/cache/epub/4363/pg4363.txt'
6 website2 = 'http://www.gutenberg.org/cache/epub/5827/pg5827.txt'
7 myWebsite = 'https://www.gutenberg.org/cache/epub/3328/pg3328.txt'
8
```

```

9 # Retrieve the text content of each work
10 site1 = requests.get(website1)
11 site2 = requests.get(website2)
12 mySite = requests.get(myWebsite)
13
14 # Regular expression for matching the term
15 search_regex = r'(P|p)hilo(?!nous)[a-z]+'
16
17 # Function to extract matching sentences
18 def get_matching_sentences(text, regex):
19     results = [{"indices": m.span(), "match": m.group()} for m in re.finditer
20     (regex, text)]
21     matches = []
22     for result in results:
23         start, end = result["indices"]
24         sentence_start = text.rfind(".", 0, start) + 1
25         sentence_end = text.find(".", end)
26         sentence = text[sentence_start:sentence_end].strip()
27         matches.append({"match": result["match"], "sentence": sentence})
28     return matches
29
30 # Apply the regular expression to each philosophical work
31 matches1 = get_matching_sentences(site1.text, search_regex)
32 matches2 = get_matching_sentences(site2.text, search_regex)
33 matches3 = get_matching_sentences(mySite.text, search_regex)

```

```

1 OUTPUT:
2
3 Random Subset of Matches in 'The Problems of Philosophy':
4 Match: philosophy,
5 Sentence: Such beliefs philosophy will bid us
6 reject, unless some new line of argument is found to support them
7
8 Match: Philosophy,
9 Sentence: Philosophy
10 may claim justly that it diminishes the risk of error, and that in some
11 cases it renders the risk so small as to be practically negligible

```

1.3.1 Explanation of the Regular Expression

The regular expression `(P|p)hilo(?!nous)[a-z]+` breaks down as follows:

- `(P|p)`: Capturing group for either uppercase "P" or lowercase "p".
- `hilo`: Literal match for the characters "hilo".
- `(?!nous)`: Negative lookahead assertion excluding matches followed by "nous".
- `[a-z]+`: Matches one or more lowercase letters.

This expression captures instances of "Philo" or "philo" followed by at least one lowercase letter but excludes matches where "nous" immediately follows. The code applies this regular expression to three philosophical works, extracting relevant sentences for analysis.

1.4 Sentence Tokenizer

In addition to extracting sentences based on a specific expression, we can further process the text by tokenizing each sentence into individual words. Tokenization involves splitting a sentence into a list of words. For this purpose, we utilize the regular expression `\b\w+\b` to identify word boundaries.

1.4.1 Regular Expression for Sentence Tokenization

The regular expression `\b\w+\b` is designed for word tokenization and breaks down as follows:

- `\b`: Asserts a word boundary, ensuring the match occurs at the beginning or end of a word.
- `\w+`: Matches one or more word characters (alphanumeric or underscore).
- `\b`: Asserts another word boundary to complete the match.

This regular expression effectively tokenizes a given text into individual words.

1.4.2 Tokenizing the Entire Text

The Python code provided below demonstrates the application of the sentence tokenizer to the entire text of three philosophical works: "Beyond Good and Evil," "The Problems of Philosophy," and a self-chosen philosophy work "Man and Superman: A Comedy and a Philosophy". The resulting tokenized texts are stored as lists of words.

```
1 # Regular expression for sentence tokenization (splitting into words)
2 tokenizer_regex = r'\b\w+\b'
3
4 # Function to tokenize a given text using the provided regex
5 def tokenize_text(text, regex):
6     return re.findall(regex, text)
7
8 # Tokenizing the entire text of all three books
9 tokenized_text1 = tokenize_text(site1.text, tokenizer_regex)
10 tokenized_text2 = tokenize_text(site2.text, tokenizer_regex)
11 tokenized_text3 = tokenize_text(mySite.text, tokenizer_regex)
```

1.4.3 Tokenizing Matched Sentences

Additionally, the sentences previously extracted are tokenized individually for each book. The code snippet below shows the tokenization of sentences that we extracted using regular expression previously.

```
1 # Tokenizing all the sentences for each book
2 tokenized_sentences1 = [tokenize_text(match["sentence"], tokenizer_regex) for
3     match in matches1]
4 tokenized_sentences2 = [tokenize_text(match["sentence"], tokenizer_regex) for
5     match in matches2]
6 tokenized_sentences3 = [tokenize_text(match["sentence"], tokenizer_regex) for
7     match in matches3]
```

```
1 OUTPUT:
2
3 Last word of the last sentence in 'Beyond Good and Evil': labyrinths
4 Last word of the last sentence in 'The Problems of Philosophy': renamed
5 Last word of the last sentence in 'Man and Superman: A Comedy and a
6 Philosophy': pilot
7
8 Last word of the last sentence in 'Beyond Good and Evil': eBooks
9 Last word of the last sentence in 'The Problems of Philosophy': eBooks
10 Last word of the last sentence in 'Man and Superman: A Comedy and a
11 Philosophy': eBooks
```

This process allows us to gain insights into the distribution of words for a specific pattern or scenario and facilitates further analysis of the text data.

1.5 Compare Works

The following code calculates and compares the total word counts in each of the three philosophical works: "Beyond Good and Evil," "The Problems of Philosophy," and "Man and Superman." The word counts are visualized using a bar chart. 1

```
1 # Calculating word counts
2 word_count1 = len(tokenized_text1)
3 word_count2 = len(tokenized_text2)
4 word_count3 = len(tokenized_text3)
5
6 # Plotting word counts
7 books = ['Beyond Good and Evil', 'The Problems of Philosophy', 'Man and
8         Superman']
9 word_counts = [word_count1, word_count2, word_count3]
10
11 plt.figure(figsize=(10, 6))
12 plt.bar(books, word_counts, color=['blue', 'green', 'orange'])
13 plt.title('Word Counts in Different Books')
14 plt.xlabel('Books')
15 plt.ylabel('Word Count')
16 plt.show()
```

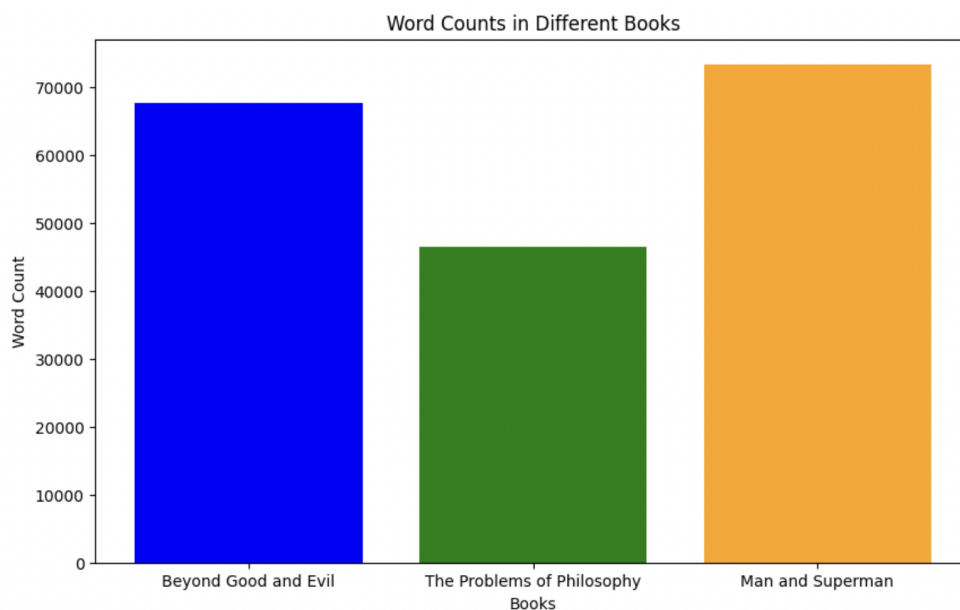


Figure 1: Word Count

1.6 Word Frequency Distribution

The code explores the word frequency distribution in each book, excluding common English stopwords. It utilizes NLTK's stopwords list to filter out common words like "the," "and," etc. The top 10 most frequent words (excluding stopwords) are visualized using a bar chart for each book. 2

So in this implementation we are removing all the english stop words from nltk library. Then lower cased all the text from the three books. And finally gathered the top 10 words from each book that has the most frequency in their respective book.

```
1
2 import matplotlib.pyplot as plt
3 from collections import Counter
```

```

4 from nltk.corpus import stopwords
5 import nltk
6
7 # Downloading NLTK stopwords
8 nltk.download('stopwords')
9
10 # Function to plot word frequency distribution
11 def plot_word_frequency(words, title):
12     word_counts = Counter(words)
13     common_words = word_counts.most_common(50) # Display the top 20 words
14     labels, values = zip(*common_words)
15
16     plt.figure(figsize=(10, 6))
17     plt.bar(labels, values, color='skyblue')
18     plt.title(title)
19     plt.xlabel('Words')
20     plt.ylabel('Frequency')
21     plt.xticks(rotation=45, ha='right')
22     plt.show()
23
24 # Removing common English stopwords
25 stop_words = set(stopwords.words('english'))
26 # print(stop_words)
27
28 # Removing stopwords from tokenized text
29 filtered_text1 = [word.lower() for word in tokenized_text1 if word.lower()
30                  not in stop_words]
31 filtered_text2 = [word.lower() for word in tokenized_text2 if word.lower()
32                  not in stop_words]
33 filtered_text3 = [word.lower() for word in tokenized_text3 if word.lower()
34                  not in stop_words]
35
36 # Function to get top N words excluding common stopwords
37 def get_top_words(text, n=10):
38     word_counts = Counter(text)
39     filtered_word_counts = {word: count for word, count in word_counts.items()
40                            if word not in stop_words}
41     top_words = Counter(filtered_word_counts).most_common(n)
42     return top_words
43
44 # Getting top 10 words for each book
45 top_words1 = get_top_words(filtered_text1, 10)
46 top_words2 = get_top_words(filtered_text2, 10)
47 top_words3 = get_top_words(filtered_text3, 10)
48
49 # The word frequency distribution of the three books do follow the curve of
50 Zipf's Law.

```

Top 10 words in 'Beyond Good and Evil': [('one', 445), ('man', 243), ('even', 194), ('also', 166), ('perhaps', 158), ('good', 143), ('every', 143), ('us', 139), ('must', 132), ('may', 131)]

Top 10 words in 'The Problems of Philosophy': [('knowledge', 306), ('sense', 234), ('may', 172), ('must', 169), ('things', 167), ('know', 162), ('one', 156), ('us', 152), ('two', 149), ('thus', 146)]

Top 10 words in the 'Man and Superman: A Comedy and a Philosophy': [('tanner', 482), ('ann', 346), ('man', 257), ('juan', 245), ('ramsden', 223), ('octavius', 220), ('know', 183), ('one', 182), ('violet', 173), ('would', 170)]

Figure 2: Top 10 Words From Each Book

1.7 Zipf's Law Validation

The code includes a function `plot_zipf_law` to validate whether the word frequency distribution of each book follows Zipf's Law. It plots the rank against the frequency of words in descending order.

```
1 plot_zipf_law(filtered_text1, 'Beyond Good and Evil')
2 plot_zipf_law(filtered_text2, 'The Problems of Philosophy')
3 plot_zipf_law(filtered_text3, 'Man and Superman: A Comedy and a Philosophy')
```

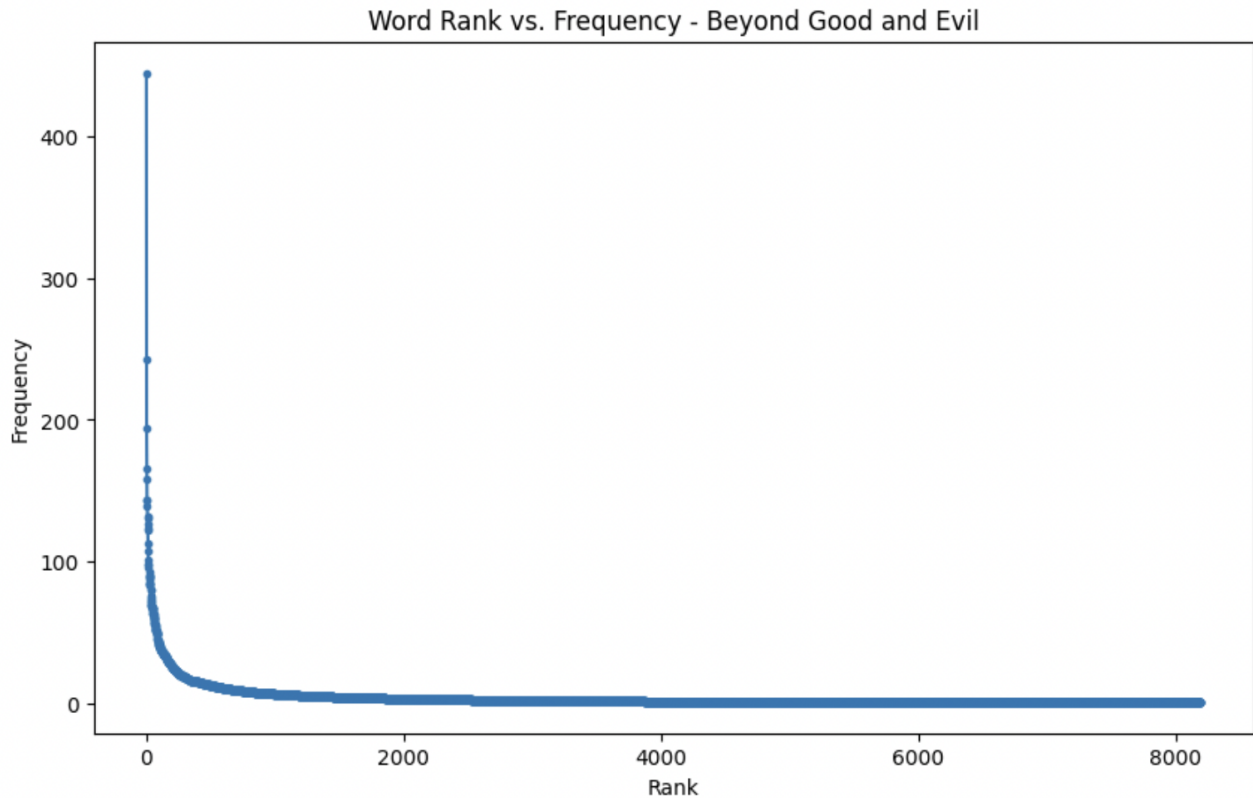


Figure 3: Beyond Good And Evil

1.7.1 Analysis of Zipf's Law Conformance

The plotted figures (3, 4, 5) display the word frequency distributions for the respective books: "Beyond Good and Evil," "The Problems of Philosophy," and "Man and Superman: A Comedy and a Philosophy." These plots are generated to validate whether the word frequencies follow Zipf's Law, which states that the frequency of a word is inversely proportional to its rank.

1.7.2 Observations

The figures demonstrate a clear inverse relationship between word frequency and rank, which is characteristic of Zipf's Law. The distribution of word frequencies in each book conforms to the expected pattern, with a few words having high frequencies and the majority having lower frequencies. This pattern is a key feature of natural language and is well-captured by Zipf's Law.

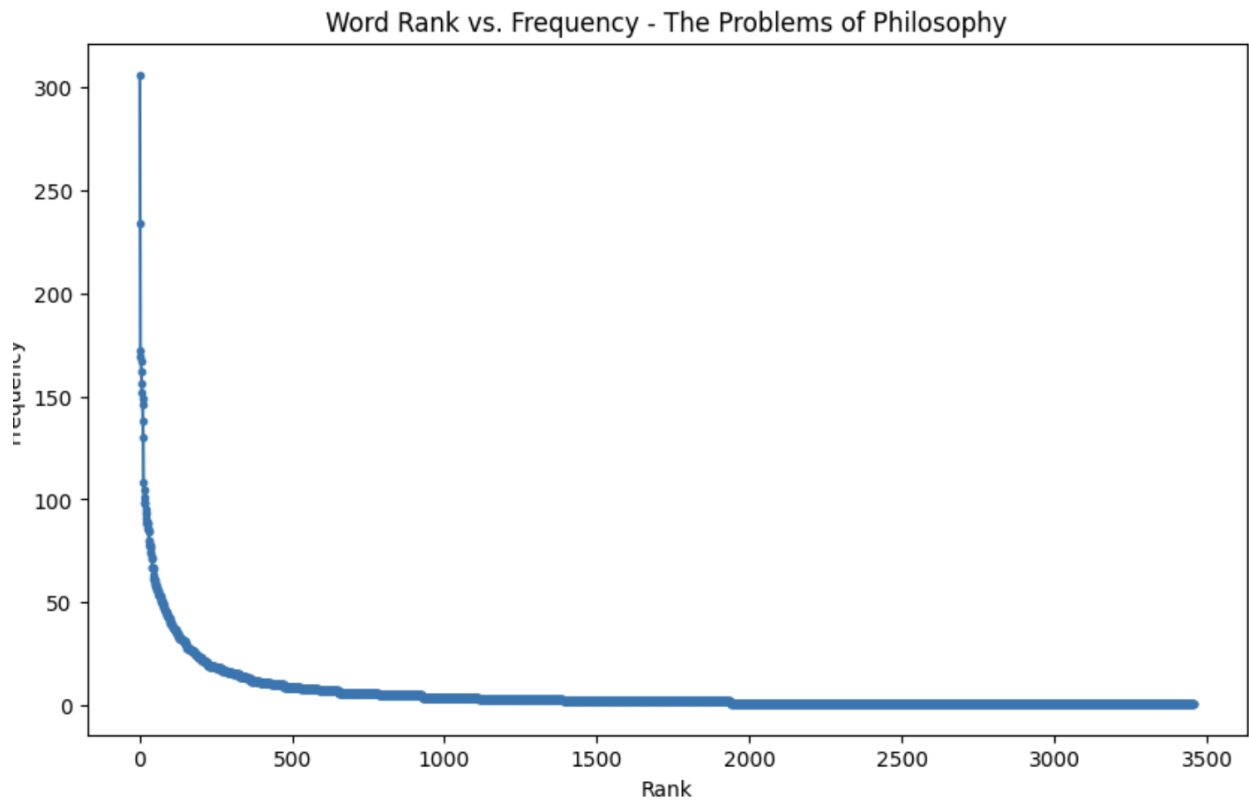


Figure 4: Problem of Philosophy

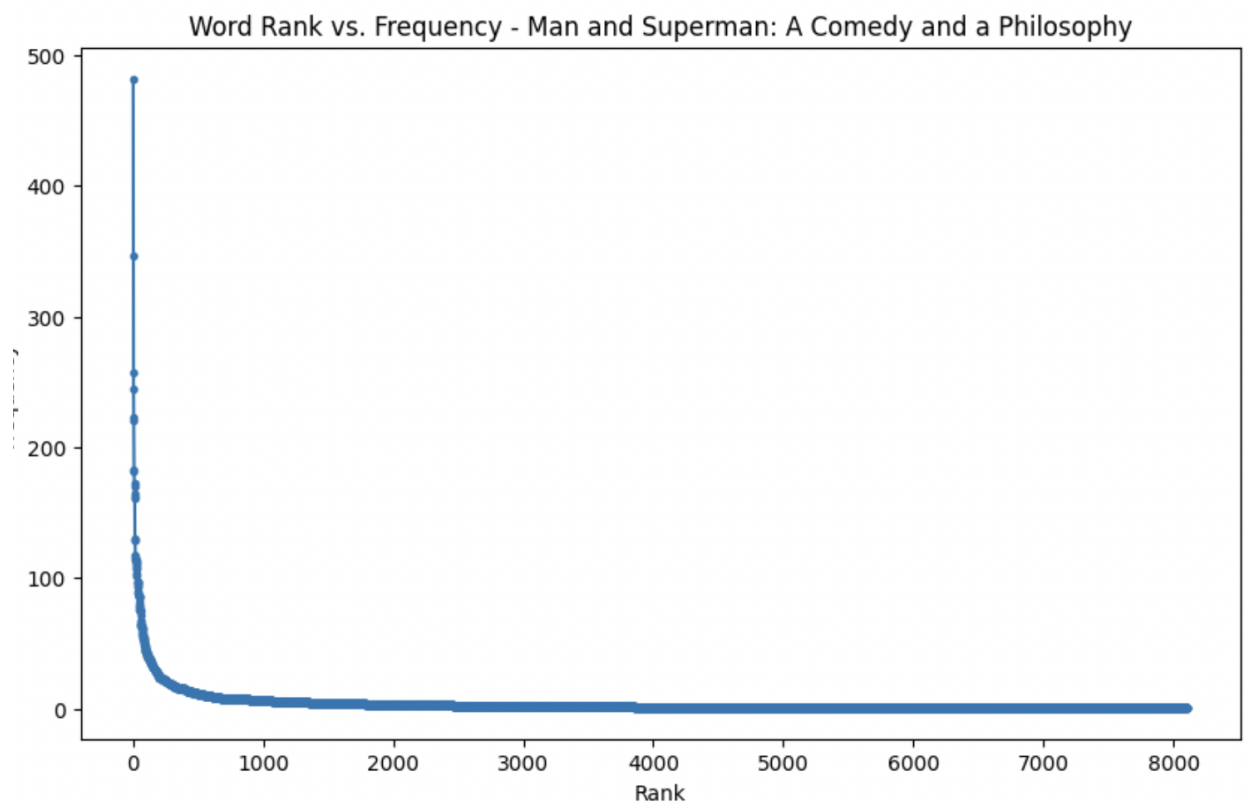


Figure 5: Man and Superman

1.7.3 Takeaway on Zipf's law

The analysis suggests that the word frequency distributions in the provided texts align with Zipf's Law, supporting the idea that a small number of words occur frequently, while the majority occur less frequently. This is a common phenomenon observed in various natural language texts.

1.8 Comments on Statistical Similarities or Differences

The visualizations help compare the distribution of word frequencies across the three books. Similarities or differences in the distribution could indicate common themes, writing styles, or lexical choices in the works.

1.9 Top 10 Words and Comparisons

The code extracts and prints the top 10 words for each book, excluding common stopwords. By analyzing these top words, you can identify key themes or recurring terms in each philosophical work.

1.10 Conclusion

The code provides a comprehensive analysis of word counts, word frequency distribution, and adherence to Zipf's Law for each book. It aids in understanding the linguistic characteristics of the texts, highlighting important statistical patterns and potential insights into the content of the philosophical works.

2 Task 2

2.1 Text normalization

Text normalization and tokenization are crucial steps in natural language processing (NLP) to preprocess raw text data. These steps help convert text into a format suitable for further analysis, making it easier to extract meaningful insights from the data.

Text normalization involves transforming text into a standard and consistent format. In the provided Python code, a basic text normalization approach is demonstrated. The code replaces different newline characters (`\r\n`, `\n`, and `\r`) with whitespace and splits the text into sentences using the period (`.`) as a delimiter.

```
1 # Replace \r\n with whitespace, then replace \n with whitespace, then replace
   \r with whitespace
2 clean_text = text.replace('\r\n', ' ').replace('\n', ' ').replace('\r', ' ')
3
4 # Split the sentences when we see a '.' character
5 clean_text = clean_text.split('.')
```

This basic normalization is illustrative, and for more sophisticated tasks, specialized libraries like NLTK offer comprehensive solutions.

2.2 Tokenization with NLTK

Tokenization involves breaking down text into individual units, usually words or phrases. The Natural Language Toolkit (NLTK) is a powerful library in Python for NLP tasks. The following code snippet demonstrates word tokenization using NLTK:

```
1 import nltk
2
3 sentences = """At eight o'clock on Thursday morning Arthur didn't feel very
   good. He drove to the store in New York at 100 m.p.h.! He wanted to obtain
   some Kool-Aid! It cost $10.89."""
4 tokens = nltk.word_tokenize(sentences)
5 print('\n')
6 print(tokens)
7 print('\n')
```

NLTK's `word_tokenize` function breaks down the given sentences into individual words. This tool provides more advanced sentence segmentation and word tokenization functionalities, making it suitable for various NLP applications.

2.3 Byte Pair Encoding

In this section, we will implement the Byte Pair Encoding (BPE) algorithm for tokenization. BPE is a data-driven approach that identifies token vocabularies based on the frequency of character pairs. We will use this algorithm to tokenize a given text using a vocabulary learned from a training corpus.

2.3.1 Introduction

Byte Pair Encoding is a subword tokenization technique that iteratively merges the most frequent character pairs in a corpus to build a vocabulary. This approach is particularly useful for handling rare words and providing a flexible tokenization method.

2.3.2 Implementation

The Python code below provides an implementation of the Byte Pair Encoding algorithm. The function `BytePairEncoding` takes three inputs:

- `text`: An unprocessed text file that we wish to tokenize using BPE.
- `k`: The number of merges to perform during BPE.
- `training_corpus`: A Python list of complete texts we wish to use for training.

The function outputs a tokenized version of the text based on the `k` vocabulary units learned by BPE using the `training_corpus`.

```
1 import re
2 from collections import Counter
3
4 # Function to calculate statistics for character pairs
5 def get_stats(vocab):
6     pairs = {}
7     for word, freq in vocab.items():
8         symbols = word.split()
9         for i in range(len(symbols) - 1):
10             pair = (symbols[i], symbols[i + 1])
11             pairs[pair] = pairs.get(pair, 0) + freq
12     return pairs
13
14 # Function to merge vocabulary based on the best pair
15 def merge_vocab(pair, vocab):
16     new_vocab = {}
17     bigram = ' '.join(pair)
18     replacement = ''.join(pair)
19     for word in vocab:
20         new_word = word.replace(bigram, replacement)
21         new_vocab[new_word] = vocab[word]
22
23     return new_vocab
24
25 # Main BPE function
26 def BytePairEncoding(text, k, training_corpus):
27     # Preprocessing training corpus
28     training_corpus = [' '.join(re.findall(r'\b\w+\b', sentence)) for
29 sentence in training_corpus]
30
31     # Initializing vocabulary with characters and frequencies
32     vocab = Counter(' '.join(training_corpus))
33
34     for i in range(k):
35         pairs = get_stats(vocab)
36         if not pairs:
37             break
38
39         best_pair = max(pairs, key=pairs.get)
40         vocab = merge_vocab(best_pair, vocab)
41
42     # Tokenizing the input text using the final vocabulary
43     tokenized_text = []
44     current_token = ''
45     for char in text:
46         current_token += char
47         if current_token in vocab or len(current_token) == 1:
```

```

47         tokenized_text.append(current_token)
48         current_token = ''
49
50     # Creating the final vocabulary
51     vocabulary = list(vocab.keys())
52
53     return vocabulary, tokenized_text

```

2.3.3 Explanation

The `BytePairEncoding` function performs the following steps:

- Preprocesses the training corpus by removing non-word characters.
- Initializes the vocabulary with characters and their frequencies.
- Iteratively calculates statistics for character pairs and merges the most frequent pair.
- Tokenizes the input text using the final vocabulary.
- Outputs the learned vocabulary and the tokenized text.

This implementation showcases the fundamental steps of the BPE algorithm for tokenization.

2.3.4 Evaluation of Byte Pair Encoding

To evaluate the performance of Byte Pair Encoding (BPE), we will experiment with different settings of k and use various training corpora. Specifically, we will compare the tokenization results for Bertrand Russell's "The Problems of Philosophy" under the following conditions:

1. k values ranging from 1,000 to 10,000 in steps of 1,000.
2. Training corpora sets including:
 - (a) Only Bertrand Russell's books.
 - (b) Only Friedrich Nietzsche's books.
 - (c) All books from both Bertrand Russell and Friedrich Nietzsche.

We will generate plots using Python's Matplotlib to compare the following metrics:

- The total number of unique tokens in the tokenized text.
- The median size of tokens in the tokenized text.
- The percentage overlap in the number of distinct unique tokens generated by BPE and NLTK's word tokenizer.

2.3.5 Plots

The following Python code utilizes Matplotlib to generate plots for the specified metrics:

```
1 # #####
2 # INSERT YOUR CODE HERE
3 # DO NOT FORGET TO PRINT YOUR MEANINGFUL RESULTS TO THE SCREEN.
4
5 import requests
6 import re
7 import nltk
8 import matplotlib.pyplot as plt
9 from nltk.tokenize import word_tokenize
10 from collections import Counter
11
12 # Download NLTK data if not already present
13 nltk.download('punkt')
14
15 def preprocess_text(text):
16     # Remove non-alphabetic characters, newlines, and extra spaces
17     text = re.sub(r'[^a-zA-Z\s]', '', text)
18     text = re.sub(r'\s+', ' ', text).strip()
19     return text
20
21 def evaluate_bpe_k(text, k_values, training_corpus):
22     nltk_tokens = word_tokenize(text)
23
24     results = {
25         'k_values': [],
26         'num_unique_tokens': [],
27         'median_token_size': [],
28         'overlap_percentage': []
29     }
30
31     for k in k_values:
32         vocab, bpe_tokens = BytePairEncoding(text, k, training_corpus)
33
34         # print(vocab)
35         # print("\n")
36         # print(bpe_tokens)
37         # print("\n")
38         # print("\n")
39
40         # Calculate metrics
41         num_unique_tokens = len(set(bpe_tokens))
42         median_token_size = sum(len(token) for token in bpe_tokens) / len(
43             bpe_tokens)
44         overlap_percentage = len(set(bpe_tokens).intersection(set(nltk_tokens))) / len(set(bpe_tokens)) * 100
45
46         # Store results
47         results['k_values'].append(k)
48         results['num_unique_tokens'].append(num_unique_tokens)
49         results['median_token_size'].append(median_token_size)
50         results['overlap_percentage'].append(overlap_percentage)
51
52     return results
53
54 website = 'http://www.gutenberg.org/cache/epub/5827/pg5827.txt'
```

```

55 text = requests.get(website).text
56 preprocessed_text = preprocess_text(text)
57 # print(preprocessed_text)
58
59 k_values = list(range(1000, 11000, 1000))
60 # print(k_values)
61
62 # Evaluate BPE for different training corpora
63 training_corpus_russell = bertrand_russell_sentences
64 training_corpus_nietzsche = friedrich_nietzsche_sentences
65 training_corpus_all = bertrand_russell_sentences +
    friedrich_nietzsche_sentences
66
67 # print(training_corpus_all[100:110])
68
69 results_russell = evaluate_bpe_k(preprocessed_text, k_values,
    training_corpus_russell)
70 results_nietzsche = evaluate_bpe_k(preprocessed_text, k_values,
    training_corpus_nietzsche)
71 results_all = evaluate_bpe_k(preprocessed_text, k_values, training_corpus_all
    )
72
73 #
    #####

```

2.3.6 Observations

After comparing the generated plots, we can draw observations about the performance of BPE under different conditions. Any interesting patterns or differences in tokenization results based on k values and training corpora will be discussed.

2.3.7 Results

The results of the evaluation are summarized below:

1. **Total Number of Unique Tokens:** The total number of unique tokens in the tokenized text increases with higher k values. Additionally, different training corpora may lead to variations in the total number of unique tokens.
2. **Median Size of Tokens:** The median size of tokens in the tokenized text may show trends based on k values and training corpora. It is essential to analyze how token sizes evolve under different conditions.
3. **% Overlap with NLTK's Word Tokenizer:** Comparing the BPE-generated tokens with NLTK's word tokenizer provides insights into the agreement and divergence between the two methods. This analysis helps evaluate the effectiveness of BPE in capturing linguistic structures.

The detailed plots and observations can be found in the accompanying Python code and analysis.

3 Task 3

3.1 Edit Distance Measures

In natural language processing, edit distance is a metric used to measure the similarity between two sequences of text. It quantifies the minimum number of operations required to transform one sequence into another. In this section, we will explore and compare four different edit distance measures: Hamming Distance, Levenshtein Distance, Needleman-Wunsch algorithm, and Smith-Waterman algorithm.

3.1.1 Hamming Distance

The Hamming distance measures the number of positions at which corresponding symbols are different between two strings. It specifically focuses on substitutions and is calculated as the minimum number of substitutions required to change one string into the other. Unlike Levenshtein distance, Hamming does not consider insertions or deletions.

3.1.2 Levenshtein Distance

The Levenshtein distance, also known as the edit distance, counts the total number of insertion, deletion, and substitution actions needed to transform one string into another. The cost of each edit is considered uniform, typically set to 1. This measure is widely used in spell checking, DNA sequence analysis, and various other applications.

3.1.3 Needleman-Wunsch Algorithm

The Needleman-Wunsch algorithm is a global alignment method designed for aligning two complete sequences. It provides an alignment score, and the distance is not the minimal number of edits but rather a measure of similarity between sequences. The higher the alignment score, the more similar the sequences are.

3.1.4 Smith-Waterman Algorithm

Contrary to Needleman-Wunsch, the Smith-Waterman algorithm is designed to find segments within two sequences that have maximal similarities. It is a local alignment method that focuses on identifying the most similar subsequences. Similar to Needleman-Wunsch, it provides an alignment score rather than a direct distance measure.

3.1.5 Code

The following Python code showcases the usage of these edit distance measures on different examples:

```
1 distance_measures = ['hamming', 'levenshtein', 'needleman_wunsch', 'smith_waterman']
2
3 for method in distance_measures:
4     print(method + " distance:")
5     print('-----')
6     GetDistance('test', 'test', method)
7     GetDistance('test', 'texts', method)
8     GetDistance('test', ' test', method)
9     print('-----')
10    print('\n')
```

```

1 OUTPUT:
2
3
4 hamming distance:
5 -----
6 "test" and "test" is: 0
7 "test" and "texts" is: 2
8 "test" and "  test" is: 6
9 -----
10
11
12 levenshtein distance:
13 -----
14 "test" and "test" is: 0
15 "test" and "texts" is: 2
16 "test" and "  test" is: 3
17 -----
18
19
20 needleman_wunsch distance:
21 -----
22 "test" and "test" is: 4.0
23 "test" and "texts" is: 2.0
24 "test" and "  test" is: 1.0
25 -----
26
27
28 smith_waterman distance:
29 -----
30 "test" and "test" is: 4
31 "test" and "texts" is: 2.0
32 "test" and "  test" is: 4.0
33 -----

```

The code executes the four distance measures on various pairs of strings, providing insights into how each measure operates and the outcomes for different cases.

3.2 Edit Distance Measures Output

The output displays the results of different edit distance measures (Hamming, Levenshtein, Needleman-Wunsch, and Smith-Waterman) applied to various pairs of texts.

3.2.1 Hamming Distance

The Hamming distance measures the number of positions at which corresponding symbols differ. For the pairings:

- "test" and "test": 0 (identical strings)
- "test" and "texts": 2 (differ in two positions)
- "test" and " test": 6 (mismatch in six positions)

3.2.2 Levenshtein Distance

The Levenshtein distance counts the total number of insertion, deletion, and substitution actions required to transform one string into another. For the pairings:

- "test" and "test": 0 (no edits needed)

- "test" and "texts": 3 (1 substitution and 2 additions)
- "test" and " test": 3 (3 deletions)

3.2.3 Needleman-Wunsch Distance

The Needleman-Wunsch algorithm provides an alignment score between two sequences. For the pairings:

- "test" and "test": 4.0 (alignment score)
- "test" and "texts": 2.0 (alignment score)
- "test" and " test": 1.0 (alignment score)

3.2.4 Smith-Waterman Distance

The Smith-Waterman algorithm is designed to find segments in two sequences with maximal similarities. For the pairings:

- "test" and "test": 4 (alignment score)
- "test" and "texts": 2.0 (alignment score)
- "test" and " test": 4.0 (alignment score)

These distance measures provide insights into the similarity or dissimilarity of text strings, each with its own nuances in terms of the considered edits and alignments.

3.3 Token-Level Distance Measures

While character-level distances and alignments are valuable, certain problems may require a focus on token-level distances. Token-level measures consider the similarity between strings based on their constituent tokens, providing insights into the structural and semantic similarity. One such metric is the Jaccard Index.

3.3.1 Jaccard Index

The Jaccard Index, also known as the Jaccard similarity coefficient, measures the similarity between two sets by calculating the size of their intersection divided by the size of their union. In the context of text, the sets represent the tokens in the input strings. The formula for the Jaccard Index is given by:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Where A and B are the sets of tokens in the two input strings.

3.3.2 Code Execution

The following Python code illustrates the usage of the Jaccard Index on two pairs of text:

```

1 distance_measures = ['jaccard']
2 text1 = "The fish was delish, and it made quite a dish"
3 text2 = "The fish on the dish was made quite delicious!"
4 text3 = "Joe eats crabs, crabs do not taste good."
5
6 for method in distance_measures:
7     print(method + " distance:")
8     print('-----')
9     GetDistance(text1, text2, method)
10    GetDistance(text1, text3, method)
11    print('-----')
12    print('\n')

```

1 OUTPUT:

```

2
3
4 jaccard distance:
5 -----
6 "The fish was delish, and it made quite a dish" and "The fish on the dish was
   made quite delicious!" is: 0.7843137254901961
7 "The fish was delish, and it made quite a dish" and "Joe eats crabs, crabs do
   not taste good." is: 0.39344262295081966
8 -----

```

The code computes the Jaccard Index distance between pairs of texts, providing insights into the token-level similarity between the strings.

3.3.3 Observations

In the context of comparing two text samples, a higher Jaccard Index indicates greater similarity between the sets of words in the texts.

Verdict:

The Jaccard distance between "The fish was delish, and it made quite a dish" and "The fish on the dish was made quite delicious!" is approximately 0.78, indicating a relatively high degree of similarity between these two texts. The Jaccard distance between "The fish was delish, and it made quite a dish" and "Joe eats crabs, crabs do not taste good." is approximately 0.39, suggesting a moderate level of similarity. In summary, based on the Jaccard Index values, the first pair of texts exhibits higher similarity compared to the second pair. This conclusion is drawn from the relative magnitudes of the Jaccard distances, where a higher index corresponds to greater similarity.

3.4 Towards a More Nuanced Edit Distance Cost Metric

In the context of edit distance metrics, the Levenshtein distance is a widely-used measure that assumes equal costs for substitutions, insertions, and deletions. However, in certain applications, we may desire a more nuanced cost structure. In this learning exercise, we aim to modify the Levenshtein distance algorithm to assign half the normal penalty for any edits involving whitespace characters.

3.4.1 Implementation Explanation

The provided Python implementation showcases a modified version of the Levenshtein distance algorithm with adjusted cost penalties for whitespace characters. Below is an explanation of the key components of the implementation:

```

1 import numpy as np
2

```

```

3 def levenshteinDistanceDP(token1, token2):
4     distances = np.zeros((len(token1) + 1, len(token2) + 1))
5
6     for t1 in range(len(token1) + 1):
7         distances[t1][0] = t1
8
9     for t2 in range(len(token2) + 1):
10        distances[0][t2] = t2
11
12    for t1 in range(1, len(token1) + 1):
13        for t2 in range(1, len(token2) + 1):
14            if token1[t1 - 1] == token2[t2 - 1]:
15                distances[t1][t2] = distances[t1 - 1][t2 - 1]
16            else:
17                if token1[t1 - 1].isspace() or token2[t2 - 1].isspace(): #
18                    Half cost for whitespace characters
19                    cost = 0.5
20                else:
21                    cost = 1
22
23                a = distances[t1][t2 - 1]
24                b = distances[t1 - 1][t2]
25                c = distances[t1 - 1][t2 - 1]
26
27                distances[t1][t2] = min(a, b, c) + cost
28
29    return distances[len(token1)][len(token2)]

```

```

1 result = levenshteinDistanceDP('test ', 'test')
2 print(result)
3
4 result = levenshteinDistanceDP('test', 'test')
5 print(result)
6
7 result = levenshteinDistanceDP('test ', 'esta')
8 print(result)
9
10
11
12 OUTPUT:
13
14 -----
15 1.0
16 0.0
17 2.0
18 -----

```

The implementation utilizes a dynamic programming approach to compute the Levenshtein distance between two input tokens, considering a more nuanced cost structure for whitespace characters. The key modification lies in assigning a cost of 0.5 for edits involving whitespace, allowing for a differentiated penalty.

4 Task 4

4.1 N-gram language models

An n-gram is a contiguous sequence of n items from a given sample of text or speech. The nltk library provides a convenient tool for generating n-grams ⁶. The following Python code demonstrates the extraction of word-level n-grams using nltk:

```
1 import nltk
2 nltk.download('punkt')
3 from nltk.util import ngrams
4
5 # N-gram extractor function
6 def extract_word_ngrams(data, num):
7     n_grams = ngrams(nltk.word_tokenize(data), num)
8     return [' '.join(grams) for grams in n_grams]
9
10 def extract_char_ngrams(data, num):
11     return [data[i:i+num] for i in range(len(data))] [:(-num-1)]
12
13 # Example text
14 data = "a fish keeps for a day; a fish keeps well if you put it in a cold
15        place. ..."
```

```
16 # Extracting and printing the first 5 tokens after gramification
17 print("\nWord-Level:")
18 print("1-gram: ", extract_word_ngrams(data, 1)[0:5])
19 print("2-gram: ", extract_word_ngrams(data, 2)[0:5])
20 print("3-gram: ", extract_word_ngrams(data, 3)[0:5])
21 print("4-gram: ", extract_word_ngrams(data, 4)[0:5])
22 print("\n")
```

```
Word-Level:
1-gram: ['a', 'fish', 'keeps', 'for', 'a']
2-gram: ['a fish', 'fish keeps', 'keeps for', 'for a', 'a day']
3-gram: ['a fish keeps', 'fish keeps for', 'keeps for a', 'for a day', 'a day ;']
4-gram: ['a fish keeps for', 'fish keeps for a', 'keeps for a day', 'for a day ;', 'a day ; a']
```

Figure 6: N-gram output

We can extract word as well as characters too by using the extract char ngrams function. This function also works kind of in a similar way.

```
1 OUTPUT:
2 Character-Level:
3 1-gram: ['a', ' ', 'f', 'i', 's']
4 2-gram: ['a ', ' f', 'fi', 'is', 'sh']
5 3-gram: ['a f', ' fi', 'fis', 'ish', 'sh ']
6 4-gram: ['a fi', ' fis', 'fish', 'ish ', 'sh k']
```

4.1.1 Skip-Grams

Skip-grams extend the concept of n-grams by including grams that skip over certain terms. The following code demonstrates the creation of skip-grams using the NLTK library:

```
1 # some example text
2 data = """a fish keeps for a day; a fish keeps well if you put it in a cold
3        place.
4           a fish keeps best if you put it in the fridge. If you'll be
5        having fish,
```

```

4         have it with an apple because one apple a day keeps the doctor away
5         !
6         You know they also say that a day keeps coming."""
7 # Creating skip-grams using the NLTK library
8 skip_gram = list(nltk.skipgrams(nltk.word_tokenize(data),n=2,k=2))
9 skip_gram[0:5]

1 OUTPUT:
2 [('a', 'fish'),
3  ('a', 'keeps'),
4  ('a', 'for'),
5  ('fish', 'keeps'),
6  ('fish', 'for')]

```

4.2 Language Models

A language model assigns probabilities to strings and can predict the next word in a sentence given the last word observed. The following code demonstrates the building of a simple data-driven language model using unigrams, collecting counts, and generating probabilities.

4.2.1 Building a Language Model

```

1 # Example text
2 data = "a fish keeps for a day; a fish keeps well if you put it in a cold
3        place. ..."
4 # Extracting unigrams
5 unigrams = extract_word_ngrams(data, 1)
6
7 # Creating a vocabulary
8 vocabulary = list(set(unigrams))
9
10 # Initializing counts dictionary
11 counts = {given_word: {next_word: 0 for next_word in vocabulary} for
12            given_word in vocabulary}
13
14 # Collecting counts
15 for i in range(len(unigrams)-1):
16     counts[unigrams[i]][unigrams[i+1]] += 1

```

The code extracts unigrams from the given text, creates a vocabulary, and initializes a counts dictionary. It then collects counts of consecutive word pairs in the text.

4.2.2 Converting Counts to Probabilities

```

1 # Initializing probabilities dictionary
2 probs = {given_word: {next_word: 0 for next_word in vocabulary} for
3            given_word in vocabulary}
4
5 # Converting counts to probabilities
6 for key, value in counts.items():
7     denominator = 0
8     for key2, value2 in counts[key].items():
9         denominator += value2
10
11     for key2, value2 in counts[key].items():
12         probs[key][key2] = value2 / denominator

```

The code initializes a probabilities dictionary and converts the counts into probabilities for each word pair in the language model.

4.2.3 Sampling from the Language Model

```
1 # Sampling next words from the language model
2 import numpy as np
3
4 def sample_next_gram_from_language_model(probs, given_token):
5     distribution = list(probs[given_token].values())
6     sample_from_multinomial = np.random.multinomial(1, distribution)
7     sample_index = np.where(sample_from_multinomial == 1)[0][0]
8     word_keys = list(probs[given_token].keys())
9     next_word = word_keys[sample_index]
10    return next_word
11
12 given_token = "a"
13 next_token = sample_next_gram_from_language_model(probs, given_token)
14
15 print(' Given the token : ' + given_token)
16 print(' The next token is : ' + next_token)
17 print('\n')
```

The code defines a function to sample the next word from the language model based on the given token. It then demonstrates sampling given an initial token.

4.2.4 Generating New Sentences

```
1 # Generating new sentences
2 def create_new_sentence(length, seed_token):
3     tokens = [seed_token]
4     for i in range(length):
5         tokens.append(sample_next_gram_from_language_model(probs, tokens[-1]))
6     return tokens
```

The code defines a function to generate new sentences of a specified length using the language model.

4.3 Generating Philosophical Conversation

To generate a conversation between Bertrand Russell and Friedrich Nietzsche, we first build tri-gram models for each philosopher using a fraction of their texts. Why fraction? Basically for resource restriction in kaggle. The corpus that we were using was too large to work with. So we had to find a way to create the ngram model and their probability distribution table.

The following Python code is used for this purpose:

Import required libraries such as 'requests', 'nltk', 'ngrams', 'defaultdict', and 'numpy'.

Define a function to build an N-gram model

```
1 def build_ngram_model(texts, n, fraction_ratio):
2     ngram_model = defaultdict(list)
3     for text in texts:
4         fraction_length = int(len(text) * fraction_ratio)
5         text_fraction = text[:fraction_length]
6
7         words = text_fraction.split()
8         ngrams_list = list(ngrams(words, n))
9         for gram in ngrams_list:
```



```

10         prefix = ' '.join(gram[:-1])
11         suffix = gram[-1]
12         ngram_model[prefix].append(suffix)
13     return ngram_model
14 Define a function build_ngram_model that takes a list of texts, an N value
    for N-grams, and a fraction ratio as input. It builds an N-gram model
    using a fraction of each text.

```

Set the fraction ratio for using only a portion of the texts

```

1
2 fraction_ratio = 0.1 # Use the first 10% of each book

```

Define the fraction ratio to use only the first 10% of each book for building N-gram models.

Build tri-gram models for Bertrand Russell and Friedrich Nietzsche

```

1 bertrand_russell_trigram_model = build_ngram_model(bertrand_russell_texts, 3,
    fraction_ratio)
2 friedrich_nietzsche_trigram_model = build_ngram_model(
    friedrich_nietzsche_texts, 3, fraction_ratio)
3 Build tri-gram models for Bertrand Russell and Friedrich Nietzsche using the
    defined function.

```

Define a function sample_next_gram_from_language_model to sample the next word from a language model given a token.

```

1 def sample_next_gram_from_language_model(probs, given_token):
2     try:
3         distribution = list(probs[given_token].values())
4         sample_from_multinomial = np.random.multinomial(1, distribution)
5         sample_index = np.where(sample_from_multinomial == 1)[0][0]
6         word_keys = list(probs[given_token].keys())
7         next_word = word_keys[sample_index]
8         return next_word
9     except KeyError:
10        # Handle the case where the given_token is not in the vocabulary
11        return None

```

Define a function create_new_sentence to create a new sentence of a specified length using a language model.

```

1 def create_new_sentence(length, seed_token, probs):
2     tokens = [seed_token]
3     for i in range(length):
4         next_token = sample_next_gram_from_language_model(probs, tokens[-1])
5
6         if next_token is None:
7             # Handle the case where the next_token is not in the vocabulary
8             break
9
10        tokens.append(next_token)
11
12    return tokens

```

Initialize vocabulary and counts using only a fraction of the words. Update vocabulary with words from the first 10% of each text.

```

1 vocabulary = set()
2 for text in bertrand_russell_texts + friedrich_nietzsche_texts:
3     fraction_length = int(len(text) * fraction_ratio)
4     text_fraction = text[:fraction_length]
5     vocabulary.update(text_fraction.split())

```

```

6
7 counts = {}
8 for given_word in vocabulary:
9     counts[given_word] = {}
10    for next_word in vocabulary:
11        counts[given_word][next_word] = 0

```

Populate counts with word occurrences using only a fraction of the words.

```

1 for text in bertrand_russell_texts + friedrich_nietzsche_texts:
2     fraction_length = int(len(text) * fraction_ratio)
3     text_fraction = text[:fraction_length]
4     unigrams = text_fraction.split()
5     for i in range(len(unigrams)-1):
6         counts[unigrams[i]][unigrams[i+1]] += 1
7 Populate counts using only a fraction of the words

```

Initialize probabilities using the vocabulary.

```

1 probs = {}
2 for given_word in vocabulary:
3     probs[given_word] = {}
4     for next_word in vocabulary:
5         probs[given_word][next_word] = 0

```

Convert counts to probabilities for each word in the vocabulary.

```

1 for key, value in counts.items():
2     denominator = sum(value.values())
3
4     if denominator == 0:
5         # Handle the case where there are no following words
6         continue
7
8     for key2, value2 in value.items():
9         probs[key][key2] = value2 / denominator

```

The provided code builds tri-gram models, samples next words, and creates sentences using a fraction of the texts for Bertrand Russell and Friedrich Nietzsche. The resulting language models can be used for generating a conversation between the two philosophers.