Natural Language Processing

CSE-3201

Project Round 1 Report

Team: Text Techies

Submitted by:

Piyush Sharma(21ucs150)

Shreshtha Gupta (21ucs196)

Shashwat Dwivedi (21UCS190)

Abhay Gupta (21UCS002)

Book: "Inferno" By Dan Brown

GitHub Link:

https://github.com/piyushsharma20/NLP_Project-.git

INDEX

1.	Objective	3
2.	Data Description	3
3.	Importing Libraries	3
4.	Importing Book	4
5.	Text Pre Processing	4
6.	Tokenizing Our Dataset	5
7.	Calculating The Frequency Of Tokens	6
8.	Creating Word Cloud	7
9.	POS Tagging	8
10.	Extracting The Largest Chapter and Tokenization	10
11.	Bigram for extracted largest Chapter	10
12.	Shannon Game	12
13.	Fill in the blanks game with bi-gram probabilities	13
14.	Output	14
15.	Accuracy Analysis	14

Objective

The goal of this text analysis report is to thoroughly examine a provided text document in.txt format ('Inferno' by Dan Brown).

Text preprocessing is used to clean and prepare the text, tokenization is used to split it down into individual words, and stop words are removed to focus on significant material. To discover common words, a frequency distribution analysis is performed, and a word cloud is generated for visualization. The distribution of tags is studied after doing part-of-speech tagging using a Penn Treebank tag set.

In addition, a bi-gram probability table with stop words is constructed for the book's major chapter. The "Shannon game" is played in a separate chapter by filling in blanks with bi-gram probabilities and accuracy.

Data Description

The book used for this project is 'Inferno' by Dan Brown (in .txt Format).

Title: Inferno

Author: Dan Brown

Number of words: 183317

Importing Libraries

To fulfill our goals, we must import necessary libraries into the code and perform essential operations on our textbook.

The required Libraries are:

re, nltk, random, wordcloud, matplotlib, stopwords, bigrams, FreqDist, word_tokenize, ConditionalFreqDist

```
In [3]: import re
import nltk|
import random
from nltk.corpus import stopwords
from nltk import bigrams, FreqDist, word_tokenize
from nltk.probability import ConditionalFreqDist
from wordcloud import WordCloud
import matplotlib.pyplot as plt
```

Importing Book

We import our book which is in .txt format.

```
In [8]: with open('output.txt', 'r', encoding='utf8') as file:
    book_text = file.read()
```

Text Pre-Processing

Text pre-processing is an important stage in the data preparation pipeline, as it improves the quality and usability of textual data. We developed many essential pre-processing stages in the framework of our project to ensure the text is clean and ready for analysis.

One of the most important steps entailed removing running section or chapter labels, which frequently crowded the material and could potentially contribute to noise in our study. To accomplish this, we used regular expressions to discover and delete recurrent patterns. We also deleted unnecessary components like photos and tables because they didn't contribute to the textual analysis we were doing.

We aimed to build a refined text corpus that is easier to deal with through these pre-processing techniques.

```
In [9]: # Remove chapter names, pictures, tables etc.
        cleaned text = re.sub(r'\b(Chapter [IVXLCDM]+|Chapter [0-9]+|Figure [0-9]+)\b', '', book text)
        cleaned text
        k, and\ncommitment.\nFor their impressive management of the London and Milan translation sites, Leon\nRomero-Montalvo and Luc
        iano Guglielmi.\nThe bright Dr. Marta Alvarez González for spending so much time with us in Florence\nand for bringing to lif
        e the city's art and architecture.\nThe peerless Maurizio Pimponi for all he did to enhance our visit to Italy.\nAll the hist
        orians, guides, and specialists who generously spent time with me in\nFlorence and Venice, sharing their expertise: Giovanna
        Rao and Eugenia Antonucci at the\nBiblioteca Medicea Laurenziana, Serena Pini and staff at the Palazzo Vecchio; Giovanna\nGiu
        sti at the Uffizi Gallery; Barbara Fedeli at the Baptistery and Il Duomo; Ettore Vito and\nMassimo Bisson at St. Mark's Basil
        ica; Giorgio Tagliaferro at the Doge's Palace; Isabella di∖nLenardo, Elizabeth Carroll Consavari, and Elena Svalduz throughou
        t all of Venice; Annalisa\nBruni and staff at the Biblioteca Nazionale Marciana; and to the many others whom I'vefailed to me
        ntion in this abbreviated list, my sincere thanks.\nRachael Dillon Fried and Stephanie Delman at Sanford J. Greenburger Assoc
        iates for\neverything they do both here and abroad.\nThe exceptional minds of Dr. George Abraham, Dr. John Treanor, and Dr. B
        ob Helm for\ntheir scientific expertise.\nMy early readers, who provided perspective along the way: Greg Brown, Dick and\nCon
        nie Brown, Rebecca Kaufman, Jerry and Olivia Kaufman, and John Chaffee.\nThe web-savvy Alex Cannon, who, along with the team
        at Sanborn Media Factory, \nkeeps things humming in the online world.\nJudd and Kathy Gregg for providing me quiet sanctuary w
        ithin Green Gables as I wrote\nthe final chapters of this book.\nThe superb online resources of the Princeton Dante Project,
```

Digital Dante at Columbia\nUniversity, and the World of Dante.The darkest places in hell are reserved for those who maintain their neutrality\nin times of moral crisis. \nFACT:\nAll artwork, literature, science, and historical references in this nove l are real.\n"The Consortium" is a private organization with offices in seven countries. Its\nname has been changed for considerations of security and privacy.\nInferno\n is the underworld as described in Dante Alighieri's epic poem \nThe\nDivine Com edy\n, which portrays hell as an elaborately structured realm\npopulated by entities known as "shades"-bodiless souls trapped between life\nand death.\nPROLOGUE\nI AM THE\n \nShade.\nThrough the dolent city, I flee.\nThrough the eternal woe, I take fl

Tokenizing Our Dataset

- Stopword Removal: We first created a set of English stopwords using stopwords.word() function of the NLTK library.
- ➤ **Text Tokenization**: We tokenized the cleaned text, breaking it into individual words using *word_tokenize()* function.
- ➤ **Lowercasing**: To ensure consistency, we converted all tokens to lowercase using *cleaned_text.lower()* function.
- > Filtering Non-Alphabetic Words: We removed tokens that contained non-alphabetic characters, retaining only words.
- ➤ **Eliminating Stopwords**: We filtered out stopwords, ensuring that the resulting 'filtered_tokens' list contained only meaningful and relevant words for subsequent text analysis.

```
# Tokenize T and remove stop words
stop words = set(stopwords.words('english'))
tokens = word_tokenize(cleaned_text.lower())
filtered tokens = [word for word in tokens if word.isalpha() and word not in stop words]
filtered_tokens
['book',
 'seek',
 'ye',
 'shall',
 'words',
 'echoing',
 'head',
 'eminent',
 'harvard',
 'symbologist',
 'robert',
'langdon',
 'awakes',
 'hospital',
 'bed',
 'recollection',
 'got',
 'explain',
 'origin',
```

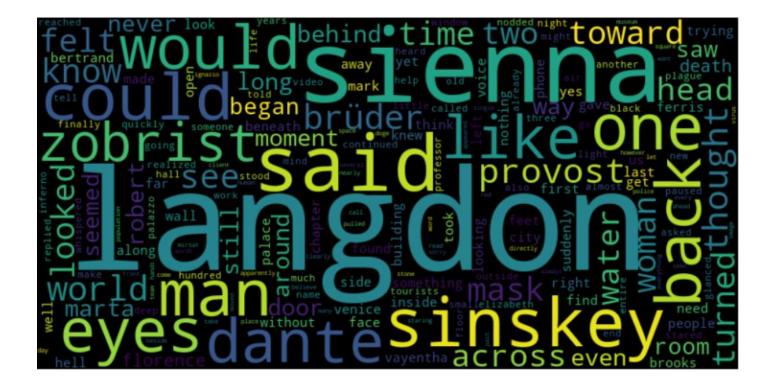
Calculating The Frequency Of Tokens

- ➤ Frequency Distribution: We used the NLTK library to create a frequency distribution of the 'filtered_tokens.' A frequency distribution counts the occurrences of each unique word in the text.
- Insights into Word Occurrence: By generating this 'token_freq_dist,' we gained valuable insights into how often different words appear in our pre-processed text data.
- ➤ **Useful for Analysis:** This frequency distribution is a foundational step for various text analysis tasks, such as identifying the most common words, understanding the vocabulary, or even spotting trends in the text corpus.

Creating Word Cloud

- ➤ Word Cloud Creation: We utilized the WordCloud library to create a visual representation of the most frequently occurring words in our text data. The 'WordCloud' object was instantiated with specific width and height parameters, defining the visual dimensions of the word cloud.
- Frequency-Based Generation: The word cloud was generated from the frequency distribution we previously created ('token_freq_dist'). This means that words that occurred more frequently in the text would be displayed more prominently in the word cloud.
- ➤ **Visualization Setup:** We configured the visualization by setting up a plot with specific dimensions (10x5 inches) to ensure clarity and readability.
- ➤ Word Cloud Rendering: The 'wordcloud' was then rendered within the plot using '*imshow*,' and 'interpolation' was set to 'bilinear' for smoother visualization.
- ➤ **Aesthetic Presentation:** We removed axis labels ('plt.axis("off")') to create a clean and aesthetically pleasing word cloud image.
- ➤ **Display:** Finally, the word cloud was displayed for visual exploration and analysis.

```
In [12]: wordcloud = WordCloud(width=800, height=400).generate_from_frequencies(token_freq_dist)
    plt.figure(figsize=(10, 5))
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis("off")
    plt.show()
```



POS Tagging

- ➤ **POS Tagging:** We applied POS tagging to our 'filtered_tokens' using NLTK (*pos_tag*) library, which assigned grammatical categories (e.g., noun, verb) to each word. Pos_tag Function of nltk library uses **Penn Treebank tagset**.
- ➤ **Tag Frequency Distribution:** The 'tag_distribution' was created to count the occurrences of different grammatical categories in the text.
- ➤ **Insights:** This analysis gives insights into the text's grammatical structure, which is valuable for understanding language usage and aids in various text analysis tasks.
- ➤ **Result Display:** We printed the most common POS tags to provide a quick overview of the distribution of grammatical categories in the text.

```
pos tags = nltk.pos tag(filtered tokens)
pos tags
[('book', 'NN'), ('seek', 'NN'),
 ('ye', 'NN'),
 ('shall', 'MD'),
 ('words', 'NNS'),
 ('echoing', 'VBG'),
 ('head', 'NN'),
 ('eminent', 'JJ'),
('harvard', 'JJ'),
 ('symbologist', 'NN'),
 ('robert', 'NN'),
 ('langdon', 'NN'),
 ('awakes', 'NNS'),
 ('hospital', 'VBP'),
 ('bed', 'JJ'),
 ('recollection', 'NN'),
 ('got', 'VBD'),
 ('explain', 'ĴJ'),
('origin', 'NN'),
```

```
tag_distribution = nltk.FreqDist(tag for word, tag in pos_tags)
tag_distribution

FreqDist({'NN': 26955, 'JJ': 14245, 'VBD': 8278, 'RB': 6032, 'NNS': 5901, 'VBG': 4644, 'VBN': 3019, 'VBP': 2392, 'IN': 1930, 'V
B': 1545, ...})

print(tag_distribution.most_common())

[('NN', 26955), ('JJ', 14245), ('VBD', 8278), ('RB', 6032), ('NNS', 5901), ('VBG', 4644), ('VBN', 3019), ('VBP', 2392), ('IN', 1930), ('VB', 1545), ('MD', 858), ('CD', 824), ('VBZ', 620), ('JJR', 270), ('RBR', 253), ('JJS', 190), ('DT', 166), ('RP', 13 0), ('FW', 127), ('PRP', 118), ('NNP', 46), ('WP$', 40), ('CC', 35), ('WDT', 33), ('UH', 11), ('RBS', 8), ('WP', 7), ('WRB', 4), ('PRP$', 1), ('POS', 1)]
```

Extracting the largest chapter and Tokenization

- ➤ Chapter Extraction: We extracted the largest chapter, namely 'Chapter C,' from our pre-processed text data. We defined the start and end indices for this chapter within the 'cleaned_text.'
- ➤ In our case, The largest chapter is Chapter 100, and Starting index = 837000, and End Index = 854415
- ➤ **Text Tokenization:** After extracting 'Chapter C,' we tokenized the text within this specific chapter, converting it to lowercase. This process prepares the text for further analysis within this particular section.
- > This action isolates 'Chapter C' for in-depth analysis or focused investigations, as needed for our project.

```
start_index_c = 837000
end_index_c = 854415
chapter_c_text = cleaned_text[start_index_c:end_index_c]
chapter_c_tokens = word_tokenize(chapter_c_text.lower())
```

Bigram for the largest chapter

- ➤ **Bigram Extraction:** We created bigrams, which are pairs of consecutive words, from the tokenized text within 'Chapter C'. This step allows us to analyze word relationships and co-occurrences in this specific section of the text.
- ➤ **Bigram Frequency Distribution:** We generated a frequency distribution ('bi_gram_freq_c') to count the occurrences of each unique bigram within 'Chapter C.' This distribution provides insights into the most common word pairs in this chapter.
- ➤ Conditional Frequency Distribution: Additionally, we calculated conditional frequencies ('bi_gram_probabilities_c') of bigrams. This

means we are not only tracking how often a bigram occurs but also how often the second word follows the first word. It can offer valuable insights into word associations and contextual relationships within 'Chapter C.'

```
bi_grams_c = list(bigrams(chapter_c_tokens))
bi_gram_freq_c = FreqDist(bi_grams_c)
bi_gram_probabilities_c = ConditionalFreqDist(bi_grams_c)
```

```
bi_grams_c

[('transpiring', 'at'),
    ('at', 'the'),
    ('the', 'cistern'),
    ('cistern', '.'),
    ('y', 'now'),
    ('now', ','),
    ('he', 'suspected'),
    ('suspected', ','),
    (',', 'sinskey'),
    ('sinskey', 'and'),
    ('and', 'the'),
    ('the', 'srs'),
    ('srs', 'team'),
    ('team', 'had'),
    ('realized', 'that'),
    ('that', 'they'),
```

```
bi_gram_freq_c
FreqDist({('.', '"'): 59, ('"', '"'): 24, (',', '"'): 23, (''', 's'): 23, ('"', 'langdon'): 22, (''', 't'): 20, ('"', 'she'): 1
8, ('"', 'i'): 15, ('.', 'he'): 14, ('of', 'the'): 14, ...})
```

➤ We will extract any other chapter (in this case consider Ch 50) and tokenize it as it is required for playing Shannon Game.

```
#Extracting any randomm chapter
start_index_other = 391725
end_index_other = 396594

chapter_other_text = cleaned_text[start_index_other:end_index_other]

chapter_other_tokens = word_tokenize(chapter_other_text.lower())
```

Shannon Game

- ➤ **Shannon Game:** We used the "Shannon Game" text modification technique within the 'chapter_other' content. The Shannon Game includes blanking out specific words in a text to see how they affect readability and comprehension. This can be a useful tool for investigating how changes in language affect comprehension.
- ➤ **Initialization:** We created an empty list called 'blanks' to hold the updated content.
- ➤ **Iteration:** We went through the 'chapter_other_tokens' to look at each word and how it related to the next word in the text.
- > Stopword Handling: If the current word is a stopword, it is simply added to the 'blanks' list. This prevents stopwords from being blanked out, which are often less informative words like "the" and "and."
- ➤ **Bigram Probability Check:** In 'bi_gram_probabilities_c,' we calculated a random probability if the current word and the following word form a known bigram. If the generated random number is smaller than the likelihood of the next word in the bigram following the current word, we added the current word to the 'blanks' list. Otherwise, we substituted "___" for a blank to simulate a word omission.
- ➤ **Else Clause:** If the current word and the following word do not create a recognized bigram, we simply add the current word to the 'blanks' list.

```
blanks = []
for i in range(len(chapter other tokens) - 1):
    current word = chapter other tokens[i]
    next word = chapter other tokens[i + 1]
    if current word in stop words:
        blanks.append(current word)
    elif (current word, next word) in bi gram probabilities c:
        if random.random() < bi gram probabilities c[current word].freq(next word):</pre>
            blanks.append(current word)
            blanks.append(" ")
    else:
        blanks.append(current word)
blanks
['et',
 'passages',
 'tour',
 'the',
 'duke',
 'of',
 'athens',
 'stairway',
 '.',
 'the',
 'sound',
 'of',
 'running',
```

Fill in the blanks game with bi-gram probabilities

- ➤ After playing the Shannon Game and filling in blanks across 'chapter_other,' we finished the chapter by adding the final word. This stage guarantees that the chapter is coherent and ends on time.
- ➤ We rebuilt the amended 'chapter_other' content by combining the components in the 'blanks' list with a space as a separator.
- Finally, we printed the filled-in chapter so that we could examine the text after using the Shannon Game. This stage allows us to examine the influence of the word changes on the readability and overall narrative flow of the text.

```
# Fill in the last word of the chapter
blanks.append(chapter_other_tokens[-1])
# Print the filled-in chapter
filled_chapter = ' '.join(blanks)
print(filled_chapter)
```

Output

et passages tour . the duke of athens stairway . the sound of running footsteps and shouting seemed to be coming from all directions now , and langdon knew their time was short . he pulled aside the curtain , and he and sienna slipped through onto a small landing . without a word , they began to descend the stone staircase . the passage had been designed as a series of frighteningly nar row switchback stairs . the deeper they went , the tighter it seemed to get . just as langdon felt as if the walls were moving in to crush him , thankfully , they could go no farther . ground le vel . the space at the bottom of the stairs was a tiny stone chamber , and although its exit had to be one of the smallest doors on earth , it was a welcome sight . only about four feet high , t he door was made of heavy wood with iron rivets and a heavy interior bolt to keep people out . " i can hear street sounds beyond the door , " sienna whispered , still looking shaken . " what 's on the other side ? " " the via della ninna , " langdon replied , picturing the crowded pedestrian walkway . " but there may be police. " " they won ' t recognize us . they ' ll be looking for a blond girl and a dark-haired man. "langdon eyed her strangely . " which is precisely what we are ... " sienna shook her head , a melancholy resolve crossing her face . " i didn 't want you to se e me like this , robert , but unfortunately it 's what i look like at the moment." abruptly , sienna reached up and grabbed a handful of her blond hair . then she yanked down , and all of her hair slid off in a single motion.langdon recoiled , startled both by the fact that sienna wore a wig and by her altered appearance without it . sienna brooks was in fact totally bald , her bare scalp smooth and pale , like a cancer patient undergoing chemotherapy . on top of it all , she 's ill ? " i know , " she said . " long story . now bend down. " she held up the wig , clearly int ending to put it on langdon 's head . is she serious ? langdon halfheartedly bent over , and sienna wedged the blond hair onto his head . the wig barely fit , but she arranged it as best as she could . then she stepped back and assessed him . not quite satisfied , she reached up , loosened his tie , and slipped the loop up onto his forehead , retightening it like a bandanna and securin g the ill- fitting wig to his head . sienna now set to work on herself , rolling up her pant legs and pushing her socks down around her ankles . when she stood up , she had a sneer on her lips . the lovely sienna brooks was now a punk-rock skinhead . the former shakespearean actress 's transformation was startling . " remember , " she said , " ninety percent of personal recognition is body language , so when you move , move like an aging rocker. " aging , i can do , langdon thought . rocker , i ' m not so sure . before langdon could argue the point , sienna had unbolted the t iny door and swung it open . she ducked low and exited onto the crowded cobblestone street . langdon followed , nearly on all fours as he emerged into the daylight . aside from a few startled gl ances at the mismatched couple emerging from the tiny door in the foundation of palazzo vecchio , nobody gave them a second look . within seconds , langdon and sienna were moving east , swallowe d up by the crowd . the man in the plume paris eyeglasses picked at his bleeding skin as he snaked through the crowd , keeping a safe distance behind robert langdon and sienna brooks . despite t heir clever disguises , he had spotted them emerging from the tiny door on the via della ninna and had immediately known who they were . he had tailed them only a few blocks before he got winded , his chest aching acutely , forcing him to take shallow breaths . he felt like he 'd been punched in the sternum . gritting his teeth against the pain , he forced his attention back to langdon and sienna as he continued to follow them through the streets of florence . chapter 50 the morning sun had fully risen now , casting long shadows down the narrow canyons that snaked between the buildings of old florence . shopkeepers had begun throwing open the metal grates that protected their shops and bars , and the air was heavy with the aromas of morning espresso and freshly baked cornetti . despite a gnawing hunger , langdon kept moving . i ' ve got to find the mask ... and see what ' s hidden on the back . as langdon led sienna northward along the slender via dei leoni , he was having a hard time getting used to the sight of her bald head . her radically altered appearance reminded him that he barely knew her . they were moving in the direction of piazza del duo mo-the square where ignazio busoni had been found dead after placing his final phone call , robert , ignazio had managed to say , breathless , what you seek is safely hidden , the gates are open to you , but you must hurry . paradise twenty-five . godspeed . paradise twenty-five , langdon repeated to himself , still puzzled that ignazio busoni had recal

Accuracy Analysis

The accuracy percentage is 21.25124131082423%. This percentage depicts the Shannon game's level of accuracy, which measures how closely the generated text matches the original text.

```
original_chapter = chapter_other_text
filled_chapter = filled_chapter

original_tokens = word_tokenize(original_chapter)
filled_tokens = word_tokenize(filled_chapter)

matching_words = sum(1 for orig, filled in zip(original_tokens, filled_tokens) if orig == filled)

accuracy_percentage = (matching_words / len(original_tokens)) * 100

print("Accuracy Percentage:", accuracy_percentage)

Accuracy Percentage: 21.25124131082423
```