# UKA TARSADIA UNIVERSITY

## BABU MADHAV INSTITUTE OF INFORMATION TECHNOLOGY

## M.SC.(IT)

## NATURAL LANGUAGE PROCESSING

## Assignment No: 1

**Student Enrolment Number : 202206100110061**

**Student Name : Divya bhaveshbhai ghori**

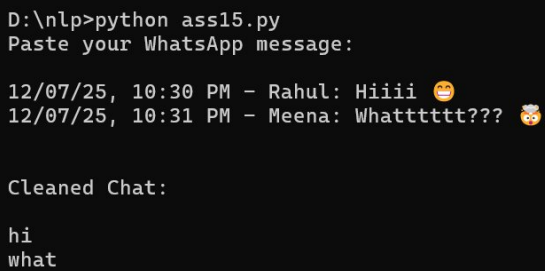| Question No: 1 | Basic Text Cleaning<br>Objective: Clean a paragraph input using Python.<br>• Take user input of a paragraph.<br>• Convert to lowercase.<br>• Remove punctuation, digits, special characters, URLs, emails, and extra spaces.<br>• Display cleaned text. |
|---|---|
| Code: | ```python<br>import re<br>import string<br><br>par = input("Enter a par: ")<br><br>par = par.translate(str.maketrans('', '', string.punctuation))<br><br>par = re.sub(r'\d+', '', par)<br><br>par = re.sub(r'[^a-zA-Z0-9\s]', '', par)<br><br>par = re.sub(r'http\S+|www\.\S+', '', par)<br><br>par = re.sub(r'\S+@\S+', '', par)<br><br>par = re.sub(r'\s+', ' ', par).strip()<br><br>print(par)<br>``` |
| Output: | ```<br>D:\nlp>python ass11.py<br>Enter a par: Uka Tarsadia University offers the BSC-IT course (3 years)! Visit: https://utu.ac.in or contact us at admission@utu.ac.in<br>.   The BMIIT department is ranked #1 in 2024!!   Apply now...<br><br>cleaned paragraph:<br>Uka Tarsadia University offers the BSCIT course years Visit or contact us at admissionutuacin The BMIIT department is ranked in Apply now<br>``` |

| Question No: 2 | Tokenization (Sentence & Word)<br>Objective: Apply sentence and word tokenization using NLTK.<br>• Input a multi-sentence paragraph.<br>• Tokenize into sentences.<br>• Tokenize each sentence into words.<br>• Display both sentence and word tokens clearly. |
| --- | --- |
| Code: | ```python<br>import nltk<br>from nltk.tokenize import sent_tokenize, word_tokenize<br><br>paragraph = input("Enter a paragraph:\n")<br><br>sentences = sent_tokenize(paragraph)<br>print("\nSentence Tokens:", sentences)<br><br>words = word_tokenize(paragraph)<br><br>print("\nWord Tokens:" , words)<br>``` |
| Output: | ```<br>D:\nlp>python ass12.py<br>Enter a paragraph:<br>Uka Tarsadia University is located in Bardoli. It offers various courses like BSC-IT and M.Sc(IT). The BMIIT department is very popular.<br><br>Sentence Tokens: ['Uka Tarsadia University is located in Bardoli.', 'It offers various courses like BSC-IT and M.Sc(IT).', 'The BMIIT department is very popular.']<br><br>Word Tokens: ['Uka', 'Tarsadia', 'University', 'is', 'located', 'in', 'Bardoli', '.', 'It', 'offers', 'various', 'courses', 'like', 'BSC-IT', 'and', 'M.Sc', '(', 'IT', ')', '.', 'The', 'BMIIT', 'department', 'is', 'very', 'popular', '.']<br>``` |

| Question No: 3 | Stop Word Removal<br>Objective: Filter meaningful words from a sentence.<br>• Input a cleaned sentence.<br>• Use NLTK to remove stop words.<br>• Print original and filtered tokens. |
| --- | --- |
| Code: | ```python<br>import nltk<br>from nltk.corpus import stopwords<br>from nltk.tokenize import word_tokenize<br><br>sentence = input("Enter a cleaned sentence:\n")<br><br>wordtokens = word_tokenize(sentence)<br>stopwords = set(stopwords.words('english'))<br>filteredtokens = [word for word in wordtokens if word not in stopwords]<br><br>print("\n Original Tokens:")<br>print(wordtokens)<br><br>print("\n Filtered Tokens :")<br>print(filteredtokens)<br>``` |

| | |
|---|---|
| Output: | ```
D:\nlp>python ass13.py
Enter a cleaned sentence:
Uka Tarsadia University is located in Gujarat and it offers various courses.

 Original Tokens:
['Uka', 'Tarsadia', 'University', 'is', 'located', 'in', 'Gujarat', 'and', 'it', 'offers', 'various', 'courses', '.']

 Filtered Tokens :
['Uka', 'Tarsadia', 'University', 'located', 'Gujarat', 'offers', 'various', 'courses', '.']
``` |

| | |
|---|---|
| Question No: 4 | Word Frequency Count<br>Objective: Analyze word usage in text.<br>• Take input of a paragraph.<br>• Clean and tokenize the text.<br>• Remove stop words.<br>• Count and display frequency of each word using collections.Counter. |
| Code: | ```
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from collections import Counter
import string

paragraph = input("Enter a paragraph:\n")

tokens = word_tokenize(paragraph)
tokens = [word.lower() for word in tokens if word.isalnum()]

stop_words = set(stopwords.words('english'))
filtered_tokens = [word for word in tokens if word not in stop_words]

word_freq = Counter(filtered_tokens)

print("\nWord Frequency Count:")

for word in word_freq:
    print(word, ":", word_freq[word])
``` |
| Output: | ```
D:\nlp>python ass14.py
Enter a paragraph:
This is a simple simple test.

Word Frequency Count:
simple : 2
test : 1
``` |

| | |
|---|---|
| Question No: 5 | You are given a copied text from a WhatsApp chat. The text contains emojis, timestamps, sender names, repeated characters like "gooooood", and random special characters.<br>Instructions:<br>Clean the chat to extract only useful message text. |

| | |
|---|---|
| | Remove timestamps, sender names, emojis, repeated characters (e.g., reduce "sooooo" to "soo"), special symbols.<br>Display cleaned messages line by line. |
| Code: | ```python<br>import re<br><br>print("Paste your WhatsApp message:\n")<br>lines = []<br>while True:<br>    line = input()<br>    if line == "":<br>        break<br>    lines.append(line)<br><br>print("\nCleaned Chat:\n")<br>for line in lines:<br>    line = re.sub(r"\d+/\d+/\d+, \d+:\d+ [AP]M - .*?:", "", line)<br>    line = re.sub(r"[^\x00-\x7F]+", " ", line)<br>    line = re.sub(r"[^a-zA-Z\s]", " ", line)<br>    line = re.sub(r"\s+", " ", line)<br>    line = re.sub(r"(.)\1{1,}", r"\1", line)<br>    line = line.lower().strip()<br>    if line:<br>        print(line)<br>``` |
| Output: | ```<br>D:\nlp>python ass15.py<br>Paste your WhatsApp message:<br><br>12/07/25, 10:30 PM - Rahul: Hiiii 😁<br>12/07/25, 10:31 PM - Meena: Whatttttt??? 🙀<br><br><br>Cleaned Chat:<br><br>hi<br>what<br>``` |

<br>

| | |
|---|---|
| Question No: 6 | Find the Most Meaningful Word<br>• Ask user to input any paragraph or story (min 5 lines).<br>• Clean and preprocess the text.<br>• Remove stop words and punctuation.<br>• Count word frequencies.<br>• Identify and print the top 3 most frequent meaningful words.<br>Input paragraph:<br>"The weather today is amazing. I love the amazing breeze and the fresh morning sunshine. The sunshine is warm and peaceful."<br>Expected Output<br>Most frequent meaningful words:<br>1. sunshine - 2 times<br>2. amazing - 2 times<br>3. breeze - 1 time |
| Code: | from collections import Counter |

| | |
|---|---|
| | ```
import string

while True:
    print("Enter at least 5 lines:")

    text = ''
    count = 0

    while True:
        line = input()
        if line == '':
            break
        text += line + ' '
        count += 1

    if count < 5:
        print("\nYou entered less than 5 lines. Enter again.\n")
        continue
    else:
        break

text = text.lower()
for ch in string.punctuation:
    text = text.replace(ch, '')

words = text.split()
word_counts = Counter(words).most_common(3)

print("\nMost frequent words:")
for i, (word, freq) in enumerate(word_counts, start=1):
    print(f"{i}. {word} - {freq} time{'s' if freq > 1 else ''}")
``` |
| Output: | ```
D:\nlp>python ass16.py
Enter at least 5 lines:
The weather today is amazing.
I love the amazing breeze.
The fresh morning sunshine.
The sunshine is warm.
It feels so peaceful today.


Most frequent words:
1. the - 4 times
2. today - 2 times
3. is - 2 times
``` |

| | |
|---|---|
| Question No: 7 | Create Your Own Text Cleaning Function<br>Students must create a function clean_text(text) that takes raw text and:<br>• Removes emails, URLs, hashtags, numbers, punctuation.<br>• Converts text to lowercase.<br>• Removes stop words.<br>Then test it on any sample review/comment/blog they find. |

| Code: | ```import re
import string

stop_words = {'a', 'an', 'the', 'is', 'in', 'on', 'and', 'or', 'at', 'to', 'for',
          'with', 'of', 'this', 'that', 'it', 'as', 'was', 'but', 'by', 'are', 'be'}

def clean_text(text):
    text = text.lower()
    text = re.sub(r'\S+@\S+', '', text)
    text = re.sub(r'http\S+|www\.\S+', '', text)
    text = re.sub(r'#\w+', '', text)
    text = re.sub(r'\d+', '', text)
    text = text.translate(str.maketrans('', '', string.punctuation))
    words = text.split()
    filtered_words = [word for word in words if word not in stop_words]
    return ' '.join(filtered_words)

inp = input("Enter review/comment/blog: ")
print("Cleaned Text:", clean_text(inp))``` |
|---|---|
| Output: | ```D:\nlp>python ass17.py
Enter review/comment/blog: Enter review/comment/blog: I loved the new movie! It's absolutely fantastic. Check it out at https://movies
.com #BestMovie
Cleaned Text: enter reviewcommentblog i loved new movie its absolutely fantastic check out``` |


| Question No: 8 | Compare Text Before and After Cleaning<br>• Input: A raw paragraph with noise (punctuation, stop words, symbols, etc.)<br>• Process: Clean it using preprocessing steps.<br>• Output:<br>o Original word count<br>o Cleaned word count<br>o Removed words list<br>o Final cleaned tokens |
|---|---|
| Code: | ```import re
import string

stop_words = {'a', 'an', 'the', 'is', 'in', 'on', 'and', 'or', 'at', 'to', 'for',
          'with', 'of', 'this', 'that', 'it', 'as', 'was', 'but', 'by', 'are', 'be'}

def clean(text):
    text = text.lower()
    text = re.sub(r'\S+@\S+', '', text)
    text = re.sub(r'http\S+|www\.\S+', '', text)
    text = re.sub(r'#\w+', '', text)
    text = re.sub(r'\d+', '', text)
    text = text.translate(str.maketrans('', '', string.punctuation))
    words = text.split()``` |

```
        cleaned = [w for w in words if w not in stop_words]
        removed = [w for w in words if w in stop_words]
        return words, cleaned, removed

text = input("Enter text: ")
original, cleaned, removed = clean(text)

print("Original Word Count:", len(original))
print("Cleaned Word Count:", len(cleaned))
print("Removed Words:", removed)
print("Cleaned Tokens:", cleaned)
```

| Output: | ``` |
|---|---|
| | D:\nlp>python ass18.py |
| | Enter text: This is a simple blog post! Visit https://example.com or email me@example.com |
| | Original Word Count: 9 |
| | Cleaned Word Count: 5 |
| | Removed Words: ['this', 'is', 'a', 'or'] |
| | Cleaned Tokens: ['simple', 'blog', 'post', 'visit', 'email'] |