



Bansilal Ramnath Agarwal Charitable Trust's  
Vishwakarma Institute of Information Technology

**Department of  
Artificial Intelligence and Data Science**

**Name:** Siddhesh Dilip Khairnar

**Class:** TY

**Division:** B

**Roll No:** 372028

**Semester:** 6<sup>th</sup>

**Academic Year:** 2023-24

**Subject Name & Code:** Natural Language and Processing & ADUA32203

**Title of Assignment:** Perform various pre-processing tasks like tokenization, stemming, lemmatization, stop word removal etc. using inbuilt functions. and using regular expressions

**Date of Performance:** 27-01-2024

**Date of Submission:** 08-02-2024

**ASSIGNMENT NO: - 2**

**Aim:** Perform various pre-processing tasks like tokenization, stemming, lemmatization, stop word removal etc. using inbuilt functions and using regular expressions

## ❖ THEORY:

- **Text Preprocessing:** Text preprocessing is a crucial step in Natural Language Processing (NLP) that involves cleaning and transforming raw text data into a format suitable for analysis. Various tasks are performed to enhance the quality of text data and improve the performance of NLP models.
- **Tokenization:**
  - **Definition:** Tokenization is the process of breaking down a text into individual units, typically words or phrases (tokens).
  - **Inbuilt Functions:** Most NLP libraries provide built-in functions for tokenization, allowing easy extraction of tokens from a given text.
  - **Regular Expressions:** Regular expressions can be employed to define custom tokenization rules, providing flexibility in handling specific patterns.
- **Stemming and Lemmatization:**
  - **Stemming:** Reducing words to their root or base form by removing suffixes.
  - **Lemmatization:** Like stemming but aims to reduce words to their canonical form (lemma), considering the context.
  - **Inbuilt Functions:** NLP libraries often offer functions for both stemming and lemmatization, helping standardize and normalize words.
  - **Regular Expressions:** Custom regular expressions can be designed to achieve stemming and lemmatization based on specific patterns.
- **Stop Word Removal:**
  - **Definition:** Stop words are common words (e.g., "the," "is," "and") that are often removed during text preprocessing to focus on content carrying words.
  - **Inbuilt Functions:** Libraries provide predefined lists of stop words and functions for their removal.

```

1  import nltk
2  from nltk.tokenize import line_tokenize, word_tokenize, TweetTokenizer
3  from nltk.stem import PorterStemmer, WordNetLemmatizer
4  from nltk.corpus import stopwords
5  from PyPDF2 import PdfReader
6
7  # Function to read PDF file and extract text
8  def read_pdf(pdf_path):
9      with open(pdf_path, 'rb') as file:
10         pdf_reader = PdfReader(file)
11         pdf_text = ''
12         for page_num in range(len(pdf_reader.pages)):
13             page = pdf_reader.pages[page_num]
14             pdf_text += page.extract_text()
15         return pdf_text
16
17  # Example PDF file path
18  pdf_path = r'D:\VS Code\Sem 6 Assignments\NLP\nlp2files\corpus.pdf'
19
20  # Read content from PDF
21  pdf_text = read_pdf(pdf_path)
22
23  # Perform line, space, word, and tweet tokenization
24  lines_tokens = line_tokenize(pdf_text)
25  space_tokens = pdf_text.split()
26  word_tokens = word_tokenize(pdf_text)
27  tweet_tokenizer = TweetTokenizer()
28  tweet_tokens = tweet_tokenizer.tokenize(pdf_text)
29
30  # Display tokenized outputs
31  print("Line Tokenization:", lines_tokens)
32  print("\nSpace Tokenization:", space_tokens)
33  print("\nWord Tokenization:", word_tokens)
34  print("\nTweet Tokenization:", tweet_tokens)
35
36  # Perform stemming using Porter Stemmer
37  porter_stemmer = PorterStemmer()
38  stemmed_tokens = [porter_stemmer.stem(token) for token in word_tokens]
39
40  # Display stemmed tokens
41  print("\nStemmed Tokens:", stemmed_tokens)
42
43  # Perform lemmatization using WordNet
44  wordnet_lemmatizer = WordNetLemmatizer()
45  lemmatized_tokens = [wordnet_lemmatizer.lemmatize(token) for token in word_tokens]
46
47  # Display lemmatized tokens
48  print("\nLemmatized Tokens:", lemmatized_tokens)
49
50  # Remove stop words using NLTK's stop-word corpus
51  stop_words = set(stopwords.words('english'))
52  filtered_tokens = [token for token in word_tokens if token.lower() not in stop_words]
53
54  # Display filtered tokens after stop-word removal
55  print("\nTokens after Stop-word Removal:", filtered_tokens)
56

```

## INPUT:

PDF

Courage is the bridge between dreams and their fulfillment in life.

## OUTPUT:

```
Microsoft Windows [Version 10.0.22631.3155]  
(c) Microsoft Corporation. All rights reserved.
```

```
D:\VS Code\Sem 6 Assignments\NLP>python -u "d:\VS Code\Sem 6 Assignments\NLP\nlp2files\NLP2pdf.py"  
Line Tokenization: ['Courage is the bridge between dreams and their fulfillment in life. ']
```

```
Space Tokenization: ['Courage', 'is', 'the', 'bridge', 'between', 'dreams', 'and', 'their', 'fulfillment', 'in', 'life. ']
```

```
Word Tokenization: ['Courage', 'is', 'the', 'bridge', 'between', 'dreams', 'and', 'their', 'fulfillment', 'in', 'life', '. ']
```

```
Tweet Tokenization: ['Courage', 'is', 'the', 'bridge', 'between', 'dreams', 'and', 'their', 'fulfillment', 'in', 'life', '. ']
```

```
Stemmed Tokens: ['courag', 'is', 'the', 'bridg', 'between', 'dream', 'and', 'their', 'fulfil', 'in', 'life', '. ']
```

```
Lemmatized Tokens: ['Courage', 'is', 'the', 'bridge', 'between', 'dream', 'and', 'their', 'fulfillment', 'in', 'life', '. ']
```

```
Tokens after Stop-word Removal: ['Courage', 'bridge', 'dreams', 'fulfillment', 'life', '. ']
```

```

1  import nltk
2  from nltk.tokenize import line_tokenize, word_tokenize, TweetTokenizer
3  from nltk.stem import PorterStemmer, WordNetLemmatizer
4  from nltk.corpus import stopwords
5
6  # Updated file path
7  corpus_path = r'D:\VS Code\Sem 6 Assignments\NLP\nlp2files\cricket_corpus.txt'
8
9  with open(corpus_path, 'r', encoding='utf-8') as file:
10     user_corpus = file.read()
11
12     # Perform line, space, word, and tweet tokenization
13     lines_tokens = line_tokenize(user_corpus)
14     space_tokens = user_corpus.split()
15     word_tokens = word_tokenize(user_corpus)
16     tweet_tokenizer = TweetTokenizer()
17     tweet_tokens = tweet_tokenizer.tokenize(user_corpus)
18
19     # Display tokenized outputs
20     print("Line Tokenization:", lines_tokens)
21     print("\nSpace Tokenization:", space_tokens)
22     print("\nWord Tokenization:", word_tokens)
23     print("\nTweet Tokenization:", tweet_tokens)
24
25     # Perform stemming using Porter Stemmer
26     porter_stemmer = PorterStemmer()
27     stemmed_tokens = [porter_stemmer.stem(token) for token in word_tokens]
28
29     # Display stemmed tokens
30     print("\nStemmed Tokens:", stemmed_tokens)
31
32     # Perform lemmatization using WordNet
33     wordnet_lemmatizer = WordNetLemmatizer()
34     lemmatized_tokens = [wordnet_lemmatizer.lemmatize(token) for token in word_tokens]
35
36     # Display lemmatized tokens
37     print("\nLemmatized Tokens:", lemmatized_tokens)
38
39     # Remove stop words using NLTK's stop-word corpus
40     stop_words = set(stopwords.words('english'))
41     filtered_tokens = [token for token in word_tokens if token.lower() not in stop_words]
42
43     # Display filtered tokens after stop-word removal
44     print("\nTokens after Stop-word Removal:", filtered_tokens)

```

---



## INPUT:

### Text:

Cricket is a popular bat-and-ball sport played between two teams, each consisting of 11 players, where the objective is to score runs by hitting the ball and running between wickets. The game is divided into innings, and teams take turns batting and bowling, with the team scoring the most runs declared the winner.

## Output:

D:\VS Code\Sem 6 Assignments\NLP>python -u "d:\VS Code\Sem 6 Assignments\NLP\nlp2files\NLP2text.py"

Line Tokenization: ['Cricket is a popular bat-and-ball sport played between two teams, each consisting of 11 players, where the objective is to score runs by hitting the ball and running between wickets. The game is divided into innings, and teams take turns batting and bowling, with the team scoring the most runs declared the winner.']

Space Tokenization: ['Cricket', 'is', 'a', 'popular', 'bat-and-ball', 'sport', 'played', 'between', 'two', 'teams', ',', 'each', 'consisting', 'of', '11', 'players', ',', 'where', 'the', 'objective', 'is', 'to', 'score', 'runs', 'by', 'hitting', 'the', 'ball', 'and', 'running', 'between', 'wickets.', 'The', 'game', 'is', 'divided', 'into', 'innings', ',', 'and', 'teams', 'take', 'turns', 'batting', 'and', 'bowling', ',', 'with', 'the', 'team', 'scoring', 'the', 'most', 'runs', 'declared', 'the', 'winner.']

Word Tokenization: ['Cricket', 'is', 'a', 'popular', 'bat-and-ball', 'sport', 'played', 'between', 'two', 'teams', ',', 'each', 'consisting', 'of', '11', 'players', ',', 'where', 'the', 'objective', 'is', 'to', 'score', 'runs', 'by', 'hitting', 'the', 'ball', 'and', 'running', 'between', 'wickets', '.', 'The', 'game', 'is', 'divided', 'into', 'innings', ',', 'and', 'teams', 'take', 'turns', 'batting', 'and', 'bowling', ',', 'with', 'the', 'team', 'scoring', 'the', 'most', 'runs', 'declared', 'the', 'winner', '.']

Tweet Tokenization: ['Cricket', 'is', 'a', 'popular', 'bat-and-ball', 'sport', 'played', 'between', 'two', 'teams', ',', 'each', 'consisting', 'of', '11', 'players', ',', 'where', 'the', 'objective', 'is', 'to', 'score', 'runs', 'by', 'hitting', 'the', 'ball', 'and', 'running', 'between', 'wickets', '.', 'The', 'game', 'is', 'divided', 'into', 'innings', ',', 'and', 'teams', 'take', 'turns', 'batting', 'and', 'bowling', ',', 'with', 'the', 'team', 'scoring', 'the', 'most', 'runs', 'declared', 'the', 'winner', '.']

Stemmed Tokens: ['cricket', 'is', 'a', 'popular', 'bat-and-bal', 'sport', 'play', 'between', 'two', 'team', ',', 'each', 'consist', 'of', '11', 'player', ',', 'where', 'the', 'object', 'is', 'to', 'score', 'run', 'by', 'hit', 'the', 'ball', 'and', 'run', 'between', 'wicket', '.', 'the', 'game', 'is', 'divid', 'into', 'inning', ',', 'and', 'team', 'take', 'turn', 'bat', 'and', 'bowl', ',', 'with', 'the', 'team', 'score', 'the', 'most', 'run', 'declar', 'the', 'winner', '.']

Lemmatized Tokens: ['Cricket', 'is', 'a', 'popular', 'bat-and-ball', 'sport', 'played', 'between', 'two', 'e', 'is', 'divided', 'into', 'inning', ',', 'and', 'team', 'take', 'turn', 'batting', 'and', 'bowling', ',', 'with', 'the', 'team', 'scoring', 'the', 'most', 'run', 'declared', 'the', 'winner', '.']

Tokens after Stop-word Removal: ['Cricket', 'popular', 'bat-and-ball', 'sport', 'played', 'two', 'teams', ',', 'consisting', '11', 'players', ',', 'objective', 'score', 'runs', 'hitting', 'ball', 'running', 'wicket', '.', 'game', 'divided', 'innings', ',', 'teams', 'take', 'turns', 'batting', 'bowling', ',', 'team', 'scoring', 'runs', 'declared', 'winner', '.']

```

1  import nltk
2  from nltk.tokenize import line_tokenize, word_tokenize, TweetTokenizer
3  from nltk.stem import PorterStemmer, WordNetLemmatizer
4  from nltk.corpus import stopwords
5  from bs4 import BeautifulSoup
6
7  # Function to read local HTML file and extract text
8  def read_local_html(file_path):
9      with open(file_path, 'r', encoding='utf-8') as file:
10         html_content = file.read()
11         soup = BeautifulSoup(html_content, 'html.parser')
12         # Extract text from HTML content
13         website_text = ' '.join([p.get_text() for p in soup.find_all('p')])
14     return website_text
15
16 # Example HTML file path
17 html_file_path = r'D:\VS Code\Sem 6 Assignments\NLP\nlp2files\corpus.html'
18
19 # Read content from the HTML file
20 website_text = read_local_html(html_file_path)
21
22 # Perform line, space, word, and tweet tokenization
23 lines_tokens = line_tokenize(website_text)
24 space_tokens = website_text.split()
25 word_tokens = word_tokenize(website_text)
26 tweet_tokenizer = TweetTokenizer()
27 tweet_tokens = tweet_tokenizer.tokenize(website_text)
28
29 # Display tokenized outputs
30 print("Line Tokenization:", lines_tokens)
31 print("\nSpace Tokenization:", space_tokens)
32 print("\nWord Tokenization:", word_tokens)
33 print("\nTweet Tokenization:", tweet_tokens)
34
35 # Perform stemming using Porter Stemmer
36 porter_stemmer = PorterStemmer()
37 stemmed_tokens = [porter_stemmer.stem(token) for token in word_tokens]
38
39 # Display stemmed tokens
40 print("\nStemmed Tokens:", stemmed_tokens)
41
42 # Perform lemmatization using WordNet
43 wordnet_lemmatizer = WordNetLemmatizer()
44 lemmatized_tokens = [wordnet_lemmatizer.lemmatize(token) for token in word_tokens]
45
46 # Display lemmatized tokens
47 print("\nLemmatized Tokens:", lemmatized_tokens)
48
49 # Remove stop words using NLTK's stop-word corpus
50 stop_words = set(stopwords.words('english'))
51 filtered_tokens = [token for token in word_tokens if token.lower() not in stop_words]
52
53 # Display filtered tokens after stop-word removal
54 print("\nTokens after Stop-word Removal:", filtered_tokens)

```

## INPUT:

### Website:

Opportunities are often disguised as hard work, so people miss them.

### Website link:

<file:///D:/VS%20Code/Sem%206%20Assignments/NLP/nlp2files/corpus.html>

## OUTPUT:

```
D:\VS Code\Sem 6 Assignments\NLP>python -u "d:\VS Code\Sem 6 Assignments\NLP\nlp2files\NLP2website.py"
```

```
Line Tokenization: ['Opportunities are often disguised as hard work, so people miss them.']
```

```
Space Tokenization: ['Opportunities', 'are', 'often', 'disguised', 'as', 'hard', 'work,', 'so', 'people', 'miss', 'them.']
```

```
Word Tokenization: ['Opportunities', 'are', 'often', 'disguised', 'as', 'hard', 'work', ',', 'so', 'people', 'miss', 'them', '.']
```

```
Tweet Tokenization: ['Opportunities', 'are', 'often', 'disguised', 'as', 'hard', 'work', ',', 'so', 'people', 'miss', 'them', '.']
```

```
Stemmed Tokens: ['opportun', 'are', 'often', 'disguis', 'as', 'hard', 'work', ',', 'so', 'peopl', 'miss', 'them', '.']
```

```
Lemmatized Tokens: ['Opportunities', 'are', 'often', 'disguised', 'a', 'hard', 'work', ',', 'so', 'people', 'miss', 'them', '.']
```

```
Tokens after Stop-word Removal: ['Opportunities', 'often', 'disguised', 'hard', 'work', ',', 'people', 'miss', '.']
```



```

1  import nltk
2  from nltk.tokenize import line_tokenize, word_tokenize, TweetTokenizer
3  from nltk.stem import PorterStemmer, WordNetLemmatizer
4  from nltk.corpus import stopwords
5  from docx import Document
6
7  # Function to read DOCX file and extract text
8  def read_docx(docx_path):
9      doc = Document(docx_path)
10     doc_text = ''
11     for paragraph in doc.paragraphs:
12         doc_text += paragraph.text + ' '
13     return doc_text
14
15 # Example DOCX file path
16 docx_path = r'D:\VS Code\Sem 6 Assignments\NLP\nlp2files\corpus.docx'
17
18 # Read content from DOCX
19 docx_text = read_docx(docx_path)
20
21 # Perform line, space, word, and tweet tokenization
22 lines_tokens = line_tokenize(docx_text)
23 space_tokens = docx_text.split()
24 word_tokens = word_tokenize(docx_text)
25 tweet_tokenizer = TweetTokenizer()
26 tweet_tokens = tweet_tokenizer.tokenize(docx_text)
27
28 # Display tokenized outputs
29 print("Line Tokenization:", lines_tokens)
30 print("\nSpace Tokenization:", space_tokens)
31 print("\nWord Tokenization:", word_tokens)
32 print("\nTweet Tokenization:", tweet_tokens)
33
34 # Perform stemming using Porter Stemmer
35 porter_stemmer = PorterStemmer()
36 stemmed_tokens = [porter_stemmer.stem(token) for token in word_tokens]
37
38 # Display stemmed tokens
39 print("\nStemmed Tokens:", stemmed_tokens)
40
41 # Perform lemmatization using WordNet
42 wordnet_lemmatizer = WordNetLemmatizer()
43 lemmatized_tokens = [wordnet_lemmatizer.lemmatize(token) for token in word_tokens]
44
45 # Display lemmatized tokens
46 print("\nLemmatized Tokens:", lemmatized_tokens)
47
48 # Remove stop words using NLTK's stop-word corpus
49 stop_words = set(stopwords.words('english'))
50 filtered_tokens = [token for token in word_tokens if token.lower() not in stop_words]
51
52 # Display filtered tokens after stop-word removal
53 print("\nTokens after Stop-word Removal:", filtered_tokens)

```

## INPUT:

Docx file:

Courage is the bridge between dreams and their realization in life.

## OUTPUT:

```
D:\VS Code\Sem 6 Assignments\NLP>python -u "d:\VS Code\Sem 6 Assignments\NLP\nlp2files\NLP2word.py"  
Line Tokenization: ['Courage is the bridge between dreams and their realization in life. ']
```

```
Space Tokenization: ['Courage', 'is', 'the', 'bridge', 'between', 'dreams', 'and', 'their', 'realization', 'in', 'life. ']
```

```
Word Tokenization: ['Courage', 'is', 'the', 'bridge', 'between', 'dreams', 'and', 'their', 'realization', 'in', 'life', '. ']
```

```
Tweet Tokenization: ['Courage', 'is', 'the', 'bridge', 'between', 'dreams', 'and', 'their', 'realization', 'in', 'life', '. ']
```

```
Stemmed Tokens: ['courag', 'is', 'the', 'bridg', 'between', 'dream', 'and', 'their', 'realiz', 'in', 'life', '. ']
```

```
Lemmatized Tokens: ['Courage', 'is', 'the', 'bridge', 'between', 'dream', 'and', 'their', 'realization', 'in', 'life', '. ']
```

```
Tokens after Stop-word Removal: ['Courage', 'bridge', 'dreams', 'realization', 'life', '. ']
```