

Ex. No. 1

07.01.26

ADDITIONAL PROGRAMS

BASIC PROGRAMS OF NLP

(A)

AIM:

: TUTORIAL

To write a python program to perform
Tokenisation by word and sentence using NLTK

Algorithm:

Step 1: Start an

Step 2: Import the required libraries

Step 3: Define a function tokensentences(text)

Step 4: Define a multi-line text input

Step 5: Call the function and store the returned
sentences in a variable.

Step 6: Iterate through the list of sentences using
a loop.

Step 7: Print each sentences with its sentence
number.

Step 8: Stop.

Program:

```
import nltk
nltk.download("punkt")
from nltk.tokenize import sent_tokenize
def tokenise_sentences(text):
    sentences = sent_tokenize(text)
    return sentences
```

OUTPUT:

(1) Profiling of morph mapping & affixes of

1) NLP is a leading platform for building
python programs to work with "human
language data".

2) It provides easy to use interface to many
50 corpora and lexical resources

(3) What are some of the things it offers?

to get started with NLTK: Python

(2) NLP is a leading platform for
building python programs to work with
"human language data".

Some of the things it offers: Python
, Redshift

.got o

: (not part)

other topics

"what how?" "background", "the
symbols, signs", "natural language", "method",
": (that) some other method of
": (that) another method of

: (that) another method of

another method

text = """NLTK is a leading platform for building python programs to work with human language data. It provides easy to use interfaces to over 50 corpora and lexical resources such as WordNet""""

Sentences = tokenise_sentences(text)
for i, sentence in enumerate(sentences):
 print(f"Sentence {i+1}: {sentence}")

(2) import nltk
nltk.download("punkt")
from nltk.tokenize import word_tokenize
def tokenize_words(text):
 words = word_tokenize(text)
 return words
text = "NLTK is a leading platform for building python programs to work with human language data"
words = tokenize_words(text)
print(words)

Result:

Thus the program to perform tokenization by word and sentences executed successfully.

(B)

AIM:

To write a python program to eliminate stopwords using NLTK

Algorithms: importing stopwords NLTK

Step 1: Start

Step 2: Import the NLTK library

Step 3: Import stopwords from nltk.corpus.

Step 4: Define the Input text

Step 5: Create a function remove_stopwords()

Step 6: Call the function with the Input text

Step 7: Store the outputs and then display it

Step 8: Stop

Program:

```
import nltk
```

```
from nltk.corpus import stopwords
```

```
from nltk.tokenize import word_tokenize
```

```
nltk.download("stopwords")
```

~~```
nltk.download("punkt")
```~~

```
def remove_stopwords(text):
```

```
 words = word_tokenize(text)
```

```
 english_stopwords = set(stopwords.words("english"))
```

All  
ability of morphing a given  
NLU given stream

Output:

NLTK leading platform building  
python programs work human language  
provide NLU w/ morphic engine

morph. HLR morph through morphic engine

best weight w/ respect to pds

(below) - encoder reading a text + 2 pds

best weight w/ HLR reading a text + 2 pds

following text best weight w/ text + 2 pds

: Pds

: Pds ; 8 pds

: morph

HLR morph

allowable morph w/ HLR easily

printed-blown beyond constraint. HLR morph

"allowable" breakdowns. HLR

("blown") breakdowns. HLR

:(text) allowable. ground. feb

(text) printed. blown : allow

(allowable) text = allowable. allow

filtered\_words = [word for word in words if  
word.lower() not in english\_stopwords]

filtered\_text = " ".join(filtered\_words)

return filtered\_text

text = "NLTK is a leading platform for building  
natural language processing programs to work with human language".  
print(filtered\_text)

filtered\_text = remove\_stopwords(text)

print(filtered\_text)

Result:

Hence the python program for removing  
stopwords executed successfully.

if `new` is not part of `word` then `word = word + new`  
else if `new` is part of `word` then `word = word - new`

Output:

"NLTK is a lead platform for building  
natural language processing tools in Python.  
It is a leading platform for research and development  
in computational linguistics and related fields like  
information retrieval, text mining, computational  
semantics, and statistical language modeling."

(c) Aim:

To write a python program to perform stemming using NLTK.

Algorithm:

Step 1: start

Step 2: Import required libraries

Step 3: Define a function stem-text(text);

Step 4: Call the function with input text

Step 5: Display the output

Step 6: Stop.

Program:

```
import NLTK nltk
```

```
def stem-text(text):
```

```
 porter_stemmer = PorterStemmer()
```

```
 words = word_tokenize(text)
```

```
 stemmed_words = [porter_stemmer.stem(word) for
 word in words]
```

```
 stemmed_text = ' '.join(stemmed_words)
```

```
 return stemmed_text.
```

text = "NLTK is a ~~leading~~ platform for building  
Python programs to work with human  
language".

```
stemmed_text = stem_text(text)
```

```
print(stemmed_text)
```

Result: Hence the python program for stemming  
executed successfully.

AIM:

To Write a python program for "Tokenization, filtration, Script validation, stopword removal, stemming the given code".

Algorithm:

Step 1: Start the program.

Step 2: Read the input text.

Step 3: Tokenize the text into words.

Step 4: Remove special characters & validate english words.

Step 5: Remove stopwords and apply stemming.

Step 6: Display the final pre-processed tokens and stop.

Program:

```
import nltk
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
nltk.download("punkt")
nltk.download("stopwords")
```

```
text = "Natural Language Processing is a fascinating
field of AI, it involves processing
human language, like English!"
```

```
def tokenize(text):
```

```
 return word_tokenize(text)
```

```
def filter_text(tokens):
```

```
 filtered_tokens = [re.sub(r'[^\w\s]', "",
 token) for token in tokens]
```

```
 filtered_tokens
```

Output:

Original: Natural Language processing  
fascinating field of AI . It  
processing like English.

Tokeized : (1) Natural Language Processing  
a fascinating field of AI  
(2) It involves processing like English

filtered token:

[ 'natural' , 'language' ] , 'processing'  
'fascinating' , 'field' , 'of' , 'AI' , 'It' , 'involves'  
'processing' , 'like' , 'English' ]

Validate Tokens: [ 'natural' , 'language' , 'process'

'is' , 'a' , 'fascinating' , 'field' , 'of' , 'A2' , '2'  
'involves' , 'processing' , 'English' ]

tokens after stopwords removal: [ 'natural' ,

'language' , 'processing' , 'fascinating' , 'field' ,  
'involves' , 'processing' , 'English' ]

Op

stemmed tokens:

[ 'Natur' , 'languag' , 'process' , 'fascin'  
'field' , 'involv' , 'process' , 'like'  
'English' ]

Result:

88

```

def validate_script(tokens):
 return [token for token in tokens if re.match(r'^[a-zA-Z]+$', token)]
def remove_stopwords(tokens):
 stop_words_set = set(stopwords.words("english"))
 return [token for token in tokens if token.lower() not in stop_words_set]
def stem_tokens(tokens):
 stemmer = PorterStemmer()
 return [stemmer.stem(token) for token in tokens]
def Preprocess_text(text):
 print("Original", text)
 print("Tokenized", tokenize(text))
 print("Filtered tokens", filter_text(tokens))
 print("Validate Tokens", validate_script(filtered_tokens))
 print("Tokens after stopwords removal", tokens_no_stopwords)
 print("Stemmed text", stemmed_tokens)
preprocessed = preprocess_text(text)

```

Result: Hence the python program for pre-processing a text executed successfully.