## LABORATORY WORK SHEET

Date: 03/02/2025

Roll No: 2275.1A.66.611. Name: B. Vaishnavi

Exp No: 02 Experiment Name: Text Preprocessing in Python

**DAY TO DAY EVALUATION:**

| | Preparation | Algorithm | Source Code | Program Execution | Viva voce | Total |
|---|---|---|---|---|---|---|
| | | Performance in the Laboratory | Calculations and Graphs | Results and Error Analysis | | |
| Max. Marks | 4 | 4 | 4 | 4 | 4 | 20 |
| Obtained | 4 | 4 | 4 | 4 | 4 | 20 |

Signature of Lab I/C

**START WRITING FROM HERE:**

2.1 Building NLP model and perform text processings prepare the text data for the NLP-model building and perform the text pre-processing. Use the required pre-processings steps based on the dataset prepared and Understand the steps involved in text pre-Processing.

Source Code:-

```
import numpy as np
import pandas as pd
import re
for df
import string
import inflect
def preprocess_text(text):
    text = text.lower()
    P = inflect.engine()
    def convert_number_to_words(match):
        return p.number_to_words(match.group())
    text = re.sub(r'\d+, convert_number_to_words, text)
    text = text.translate(str.maketrans('', '', string.punctuation)
```

```python
text = ' '.join(text.split())
return text

statements = [
    "Hello ! How are you Today?",
    "I have 2 cats and 1 dog.",
    "NLP is amazing!!!",
    "The temperature is 30 degrees celsius."
]
preprocessed_statements = [preprocess_text(sentence) for sentence
                           in statements]
for original, processed in zip(statements, preprocessed_statements):
    print(f"Original : {original} \n Preprocessed : {processed}\n")
```

Output:-

Original: Hello! How are you today?
Preprocessed: hello how are you todays
Original: I have 2 cats and 1 dog.
Preprocessed: i have two cats and one dog
Original: NLP is amazing!!
Preprocessed: nlp is amazing
original: The temperature is 30 degrees celsius
Preprocessed: the temperature is 30 degree desicy

## 2.2 TEXT PREPROCESSING OPERATIONS:

Prepare the text data for the NLP model building and perform the text pre-processing. Use the required pre-processing steps.

Source Code:-

```python
import pandas as pd
import numpy as np
import spacy
from spacy.lang.en.stop_words import STOP_WORDS as stop_words
df = pd.read_csv ('https://raw.githubuser content.com/laxmienit/csv', encoding='latin2')
df ['word_counts'] = df ['tweits'].apply(lambda x: len (str(x).split()))
df ['chor_counts'] = df ['tweits'].apply (lambda x: len (str(x)))
df ['avg_word_length'] = df ['chor_counts'] / df ['word_counts']
print ("Sample Data:")
print (df.sample (5))
print ("\nMax Word Count :", df ['word_counts'].max())
```

```python
Print ("Min Word Count :", df['word_counts'].min())
Print ("In Tweets with one word:")
print (df[df['word_counts'] == 1])
```

output:-

Max word Count : 32
Min word Count : 1

Tweets with One word :

| | twitts | Sentiment | word_Counts | Char_Counts |
|---|---|---|---|---|
| 385 | homework | 0 | 1 | 9 |
| 691 | @entelly | 0 | 1 | 9 |
| 1124 | disappointed | 0 | 1 | 13 |
| 1286 | @officedrngnfox | 0 | 1 | 16 |
| 1325 | headache | 0 | 1 | 9 |
| ... | | | | |
| 2947 | 8.0 | | | |
| 3176 | 13.0 | | | |

## 2.3 Preprocessing and cleaning.

Prepare the text data for NLP model building and perform the pre-processing.

Source Code:-

```python
import pandas as pd
import numpy as np
import spacy
import re
from spacy.lang.en.stop_words import STOP_WORDS as stopwords
from contractions import fix
df = pd.read_csv('https://raw.githubusercontent.com/laxmimerit/twitter-data/master/twitter4000.csv', encoding='latin1')
def preprocess_text(text):
    if not isinstance(text, str):
        return ""
    Text = text.lower()
    text = fix(text)
    text = re.sub(r'\s+@\s+', '', text)
    return Text
df['clean_text'] = df['twitts'].astype(str).apply(preprocess_text)
print ("Gmail count:", df['clean_text'])
df['email_count'] = df['twitts'].apply(lambda x: len(re.findall(r'\s+@\s+', str(x))))
Print ("Email Count:", df['email_count'])
```

Email Count : 0

| | |
|---|---|
| 1 | 0 |
| 2 | 0 |
| 3 | 0 |
| 4 | 0 |
| .. | .. |
| 3995 | 0 |
| 3996 | 0 |
| 3997 | 0 |
| 3998 | 0 |
| 3999 | 0 |
| ... | |

| | | |
|---|---|---|
| 2292 | in line for the simpsons slide. more cute foreign... | 0 |
| 1311 | need to save up for this Sexail $70 dress ...g... | 0 |

## 2.4 Preprocessing and cleaning.
implement the text pre procesing and perform whirous operations.

Source Code :-

```
import pandas as pd
import re unicodedata nltk
from bs4 import BeautifulSoup
from nltk. Corpus import stopwords
nltk.download ('stopwords')
print (df.columns)
def remove_html_tags (twitts):
    return BeautifulSoup (twitts, 'html.parser').get_text()
def remove_accented_chars(text):
    return ''.join (c for c in unicodedata.normalize ('NFKD', text) if not
                    unicodedata.Combining(c))
def remove_stopwords (text):
    stop_words = Set (stopwords.words ('english'))
    return ''.join (word for word in words if word.lower() not in shop_words)
df ['cleaned-text'] = df ['twitts'].astype (str)
df['cleaned-text'] = df ['cleaned-text'].apply (remove_html_tags)
df['cleaned-text'] = df ['cleaned_text'].apply (remove_accented_chars)
df['cleaned-text'] = df['cleaned-text'].apply (remove_stopwords)
df.to-csv ('twitter_4000,cleaned.csv', index = False
print ("Text preprocessing completed, cleaned data saved to 'twitter_4000_
                                                            cleaned.csv")
```

output :-

Text preprocessing Completed. Cleaned data Saved to
'twitter4000-cleaned.csv'.
[nltk_data] Downloading package stopwords to /user/share /nltk_data.
[nltk_data] package stopwords is already up_to_date!

## 2.5 PREPROCESSING and cleaning:

Prepare the text data for the NLPmodel building and perform text preprocessing

Source code:

```python
import pandas as pd
import re   unicodedata   nltk
from bs4 import BeautifulSoup
from nltk.corpus import stopwords
from textblob import TextBlob
from wordcloud import WordCloud
nltk.download('stopwords')
df = pd.read_csv('https://raw.githubusercontent.com/laxminit/twitter.csv, encoding='latin1')

def correct_spelling(text):
    return str(TextBlob(text).correct())

def tokenize_text(text):
    return TextBlob(text).words
df['cleaned_text'] = df['cleaned_text'].apply(correct_spelling)
df['tokens'] = df['cleaned_text'].apply(tokenize_text)
def generate_word_cloud(text):
    wordcloud = WordCloud(width=800, height=400, bg_clr='white').generate(' '.join(text))
    plt.figure(figsize=(10,5))
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis('off')
    plt.show()
generate_word_cloud(df['cleaned_text'])
df.to_csv('twitter4000_cleaned.csv', index=False)
print("Text Preprocessing Completed. cleaned data saved to twitter4000_cleaned.csv.")
```

Output:

```
[nltk_data] Downloading package stopwords to /usr/share/nltk_data...
[nltk_data] package stopwords is already up-to-date!
```

 400 Hundred

Image
800