

NLP EXPERIMENT NO : 2

CODE :

```
# Install required packages

!pip install nltk regex

import nltk

import re

import string

import pandas as pd

from nltk.tokenize import word_tokenize

from nltk.corpus import stopwords

import zipfile

import urllib.request

# Download NLTK resources

nltk.download('punkt')

nltk.download('stopwords')

# Download and extract dataset

dataset_url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/00228/smsspamcollection.zip'

filename = 'smsspamcollection.zip'

urllib.request.urlretrieve(dataset_url, filename)

with zipfile.ZipFile(filename, 'r') as zip_ref:

    zip_ref.extractall('sms_spam_data')

# Load dataset

data_path = 'sms_spam_data/SMSSpamCollection'

# Read dataset (tab-separated)

df = pd.read_csv(data_path, sep='\t', header=None, names=['label', 'message'])
```

```
# Show first 5 rows

print(df.head())

# Preprocessing function

def preprocess_text(text):

    # Tokenize text into words

    tokens = word_tokenize(text)

    # Remove stopwords and punctuation

    stop_words = set(stopwords.words('english'))

    filtered_tokens = [word for word in tokens if word.lower() not in stop_words and word not in
string.punctuation]

    # Keep only ASCII characters (Latin script)

    latin_tokens = [word for word in filtered_tokens if re.match(r'^[\x00-\x7F]+$', word)]

    return ' '.join(latin_tokens)

# Apply preprocessing to the messages

df['cleaned_message'] = df['message'].apply(preprocess_text)

# Display cleaned messages

df[['message', 'cleaned_message']].head()
```

OUTPUT :

cleaned_message

