**[Day-7 2211cs020196]Write a Python script that: 1. Use Genism to preprocess data from a sample text file, follow basic procedures like tokenization, stemming, lemmatization etc. ¶**

```
In [3]:   1  !pip install gensim nltk spacy
          2  import re
          3  import gensim
          4  from nltk.stem.porter import PorterStemmer
          5  from nltk.corpus import stopwords
          6  import spacy
          7  import nltk
          8  nltk.download('stopwords')
          9  nlp = spacy.load("en_core_web_sm")
         10  porter_stemmer = PorterStemmer()
         11  stop_words = set(stopwords.words('english'))
         12  def preprocess_text(text):
         13      text = re.sub(r'[^\w\s]', '', text.lower())
         14      tokens = [word for word in gensim.utils.simple_preprocess(text) if wor
         15      stemmed_tokens = [porter_stemmer.stem(token) for token in tokens]
         16      doc = nlp(' '.join(stemmed_tokens))
         17      lemmatized_tokens = [token.lemma_ for token in doc]
         18      return lemmatized_tokens
         19  text_content = """
         20  Write a Python script that uses Gensim to preprocess data from a sample te
         21  file. Follow basic procedures like tokenization, stemming, and lemmatizati
         22  Print the final output to verify the preprocessing steps.
         23  """
         24  processed_text = preprocess_text(text_content)
         25  print(processed_text)
         26
```

['write', 'python', 'script', 'use', 'gensim', 'preprocess', 'data', 'sampl',
'text', 'file', 'follow', 'basic', 'procedur', 'like', 'token', 'stem', 'lemm
at', 'print', 'final', 'output', 'verifi', 'preprocess', 'step']

In [ ]:    1