

1 **# [Day-8 2211cs020208] Write a python script that:1.Tokenizes a sample paragraph into words and sentence**

In [*]:

```
1 !pip install nltk
2 import nltk
3 nltk.download('punkt')
4 def tokenize_text(paragraph):
5     sentences = nltk.sent_tokenize(paragraph)
6     words = [nltk.word_tokenize(sentence) for sentence in sentences]
7     return sentences, words
8 paragraph = """
9 Tokenization is the process of breaking text into smaller units called to
10 These tokens can be words, sentences, or even smaller units. Tokenization
11 important step in text preprocessing.
12 """
13 sentences, words = tokenize_text(paragraph)
14 print("Sentences:")
15 for sentence in sentences:
16     print(sentence)
17 print("\nWords:")
18 for word_list in words:
19     print(word_list)
```

Requirement already satisfied: nltk in c:\users\nihar\anaconda3\lib\site-packages (3.7)

Requirement already satisfied: click in c:\users\nihar\anaconda3\lib\site-packages (from nltk) (8.0.4)

Requirement already satisfied: joblib in c:\users\nihar\anaconda3\lib\site-packages (from nltk) (1.1.0)

Requirement already satisfied: regex>=2021.8.3 in c:\users\nihar\anaconda3\lib\site-packages (from nltk) (2022.7.9)

Requirement already satisfied: tqdm in c:\users\nihar\anaconda3\lib\site-packages (from nltk) (4.64.1)

Requirement already satisfied: colorama in c:\users\nihar\anaconda3\lib\site-packages (from click->nltk) (0.4.6)

In []:

1