

Sardar Vallabhbhai National Institute of Technology (SVNIT) Surat
Department of Artificial Intelligence
B.Tech. Artificial Intelligence

B. Tech. III (AI) Semester – V NATURAL LANGUAGE PROCESSING AI357 Scheme	L	T	P	Credit
	3	0	2	04

Assignment 1 Text Preprocessing

1. You need to complete 4 tasks.
 - a. Visit <https://huggingface.co/datasets/ai4bharat/IndicCorpV2> website and download the data from your language. Extract all the data.
 - b. You need to write codes for a sentence tokenizer and word tokenizer. Tokenize each paragraph into sentences and words. Tokenize each word. Your tokenizer should tokenize punctuations, URLs, numbers (handle decimals), mail ids, dates.
 - c. After your data is tokenized, save them into a file or multiple files.
 - d. Then compute the following corpus statistics:
 - i. Total number of sentences
 - ii. Total number of words
 - iii. Total number of characters
 - iv. Average Sentence Length (Average number of words per sentence)
 - v. Average word length (Average number of characters per word)
 - vi. Type/Token Ratio (TTR) (Total number of unique tokens / Total number of tokens)
2. Repeat the same steps on a huge monolingual corpora available at <https://huggingface.co/datasets/oscar-corpus/OSCAR-2301>

Instructions:

Create a folder "NLP" in your GitHub. For every assignment, create a folder like "Assignment_1", "Assignment_2", and so on. Put your codes in the respective folders. Each folder should be associated with a complete README.md file. Your lab evaluation will be based on the codes you upload.

For the 1st assignment, you need to store the data in the below format.

Write lines to the file where you are storing the tokenized data. Each line should contain a tokenized sentence. After you tokenize each paragraph into sentences and each sentence into words, combine the tokenized words by spaces to form tokenized sentences and write them to a file.

Sardar Vallabhbhai National Institute of Technology (SVNIT) Surat

Department of Artificial Intelligence

B.Tech. Artificial Intelligence

As the file sizes are large, use compression techniques and save them as parquet (learn about parquet files).