

AIWIR Lab Week 2

Team-8

Team members :

1. A Spoorthi Alva PES2UG19CS001
2. A R Manyatha PES2UG19CS002
3. Achyut Jagini. PES2UG19CS013
4. Amulya S Dinesh PES2UG19CS035

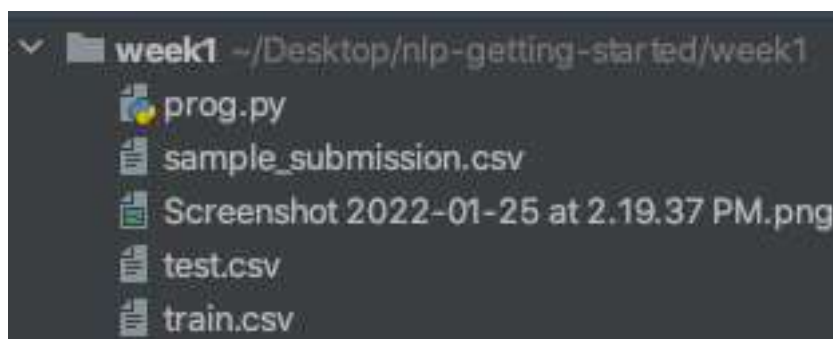
Dataset: <https://www.kaggle.com/c/nlp-getting-started/data?select=train.csv>

Tasks :

- Tokenize each Tweet into sentences
- Tokenize each tweet into words
- Remove stopwords in each tweet - NLTK library

Tool used : Pycharm

- It is installed on the computer.
- The dataset is downloaded and added to the folder where the main program exists
- Project is created in Pycharm by opening a folder containing the dataset.
- Required modules are installed.



Code :

```
import nltk
nltk.download('punkt')
nltk.download('omw-1.4')
nltk.download('stopwords')
from nltk.tokenize import word_tokenize
from nltk.stem import PorterStemmer
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
nltk.download('wordnet')
stop_words=set(stopwords.words('english'))
#using stemming-TEXT NORMALIZATION
file='train.csv'
filer=open(file,"r",encoding="utf-8")
text=filer.read()
text=text.replace("\n"," ")
#tokenizing words
word_tokens = word_tokenize(text)

filtered_sentence = []
symbols=[':','/',' ','(',')','@','?',';','//','#','!','&','$','%','*','...','.',',','..','-','[',']','{','}']
for w in word_tokens:
    if w not in stop_words and w not in symbols:
        filtered_sentence.append(w)
print("The text after stemming" )
Stem_words = []
ps = PorterStemmer()

for w in filtered_sentence:
    rootWord = ps.stem(w)
    Stem_words.append(rootWord)
print(filtered_sentence)
print(Stem_words)

lemma_word = []

wordnet_lemmatizer = WordNetLemmatizer()
for w in filtered_sentence:
    word1 = wordnet_lemmatizer.lemmatize(w, pos = "n")
    word2 = wordnet_lemmatizer.lemmatize(word1, pos = "v")
    word3 = wordnet_lemmatizer.lemmatize(word2, pos = ("a"))
    lemma_word.append(word3)
print("The text after lemmatization")
print(lemma_word)
```

