

Stop words are available in abundance in any human language. By removing these words, we remove the low-level information from our text in order to give more focus to the important information. In other words, we can say that the removal of such words does not show any negative consequences on the model we train for our task.

Removal of stop words definitely reduces the dataset size and thus reduces the training time due to the fewer number of tokens involved in the training.

```
import nltk
from nltk.corpus import stopwords
nltk.download('stopwords')

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
True
```

```
stopwords.words("english")          #stopwords.words("spanish")
```

```
['i',
 'me',
 'my',
 'myself',
 'we',
 'our',
 'ours',
 'ourselves',
 'you',
 "you're",
 "you've",
 "you'll",
 "you'd",
 'your',
 'yours',
 'yourself',
 'yourselves',
 'he',
 'him',
 'his',
 'himself',
 'she',
 "she's",
 'her',
 'hers',
 'herself',
 'it',
 "it's",
 'its',
 'itself',
 'they',
 'them',
 'their',
 'theirs',
 'themselves',
 'what',
 'which',
 'who',
 'whom',
 'this',
 'that',
 "that'll",
 'these',
 'those',
 'am',
 'is',
 'are',
 'was',
 'were',
 'be',
 'been',
 'being',
 'have',
 'has',
 'had',
 'having',
```

```
'do',  
..  
  
text = "The 2018  Fifa World Cup was held in Russia.France won the Tournament and Croatia were the semifinalist."  
  
stopwords = set(stopwords.words("english"))  
  
from nltk.tokenize import word_tokenize  
nltk.download('punkt')  
  
[nltk_data] Downloading package punkt to /root/nltk_data...  
[nltk_data]   Unzipping tokenizers/punkt.zip.  
True  
  
result = word_tokenize(text)  
  
final=[]  
  
for i in result:  
    if i not in stopwords:  
        final.append(i)  
        final.append(" ")  
  
"".join(final)  
  
'The 2018 Fifa World Cup held Russia.France Tournament Croatia semifinalist . '
```

[Colab paid products](#) - [Cancel contracts here](#)

✓ 0s completed at 12:54 PM

