

ir-01-1

October 9, 2024

```
[10]: # !pip install nltk
```

```
[26]: import nltk
      from nltk.corpus import stopwords
      nltk.download('stopwords')
      print(stopwords.words('english'))
```

```
['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', 'does', 'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or',
'because', 'as', 'until', 'while', 'of', 'at', 'by', 'for', 'with', 'about',
'against', 'between', 'into', 'through', 'during', 'before', 'after', 'above',
'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under',
'again', 'further', 'then', 'once', 'here', 'there', 'when', 'where', 'why',
'how', 'all', 'any', 'both', 'each', 'few', 'more', 'most', 'other', 'some',
'such', 'no', 'nor', 'not', 'only', 'own', 'same', 'so', 'than', 'too', 'very',
's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now',
'd', 'll', 'm', 'o', 're', 've', 'y', 'ain', 'aren', "aren't", 'couldn',
"couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', "hadn't", 'hasn',
"hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't",
'mustn', "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn',
"shouldn't", 'wasn', "wasn't", 'weren', "weren't", 'won', "won't", 'wouldn',
"wouldn't"]
```

```
[nltk_data] Downloading package stopwords to
```

```
[nltk_data] C:\Users\HP\AppData\Roaming\nltk_data...
```

```
[nltk_data] Package stopwords is already up-to-date!
```

```
[27]: import nltk
      from nltk.tokenize import word_tokenize
      from nltk.corpus import stopwords

      # Example sentence
```

```

example_sent = """This is a sample sentence, showing off the stop words_
↳filtration."""

# Set of stop words
stop_words = set(stopwords.words('english'))

# Tokenize the sentence
word_tokens = word_tokenize(example_sent)

# Converts the words in word_tokens to lower case and then checks whether they_
↳are present in stop_words or not
filtered_sentence = [w for w in word_tokens if not w.lower() in stop_words]

# Without lower case conversion
filtered_sentence = []
for w in word_tokens:
    if w not in stop_words:
        filtered_sentence.append(w)

# Print the original tokens and filtered sentence
print(word_tokens)
print(filtered_sentence)

```

```

['This', 'is', 'a', 'sample', 'sentence', ',', 'showing', 'off', 'the', 'stop',
'words', 'filtration', '.']
['This', 'sample', 'sentence', ',', 'showing', 'stop', 'words', 'filtration',
'.']

```

```

[28]: import io
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
# word_tokenize accepts
# a string as an input, not a file.
stop_words = set(stopwords.words('english'))
file1 = open("text.txt")
# Use this to read file content as a stream:
line = file1.read()
words = line.split()
for r in words:
    if not r in stop_words:
        appendFile = open('filteredtext.txt','a')
        appendFile.write(" "+r)
        appendFile.close()

```

![[image.png](attachment:87aabb3d3-2c56-4592-8863-cef4dcacf2cc.png)]![[image.png](attachment:605c211e-5057-4d8f-9b6d-06a9b0f6da5b.png)]

```
[29]: from nltk.stem import PorterStemmer
from nltk.tokenize import word_tokenize
ps = PorterStemmer()
# choose some words to be stemmed
words = ["program", "programs", "programmer", "programming", "programmers"]
for w in words:
    print(w, " : ", ps.stem(w))
```

```
program : program
programs : program
programmer : programm
programming : program
programmers : programm
```

```
[30]: from nltk.stem import PorterStemmer
from nltk.tokenize import word_tokenize
ps = PorterStemmer()
sentence = "Programmers program with programming languages"
words = word_tokenize(sentence)
for w in words:
    print(w, " : ", ps.stem(w))
```

```
Programmers : programm
program : program
with : with
programming : program
languages : languag
```

[]:

[]: