# Processing using NLTK

**Uzair Ahmad**

## ⌄ What is text preprocessing in NLP and why it is important?

**Text preprocessing is a first step in any text based classification. In large corpus, text comes from various resources which contain noise along with valueable information. If we use unprocessed data directly with the models, it model will peform badly and give unpredictable result.**

## ⌄ Steps to install NLTK

**Mac/Unix**

From the terminal:

1. Install NLTK: run `pip install -U nltk`
2. Test installation: run `python` then type `import nltk`

**Windows**

1. Install NLTK: http://pypi.python.org/pypi/nltk
2. Test installation: `Start>Python35`, then type `import nltk`

## ⌄ Download NLTK data

```
1 import nltk
2 nltk.download()
```

## ⌄ 1) Convert all words to lower case

**Lower and upper case for same words is treated as different by the models, so we need to convert all words to lower case or uppper case**

```
1 def toLowerCase(text):
2     return text.lower() #changes all upper case alphabet to lower case
```

```
1 text = 'Did you catch the bus ? Are you frying an egg ?'
2 print(f'Before - {text}')
3 print(f'After  - {toLowerCase(text)}')
4
```

```
Before - Did you catch the bus ? Are you frying an egg ?
After  - did you catch the bus ? are you frying an egg ?
```

## 2) Removal of URLs

**Regular expression or RegEx in Python is denoted as RE (REs, regexes or regex pattern) are imported through re module. The functions in this module let you check if a particular string matches a given regular expression.**

```
1 import re
2 def removeURLs(text):
3
4     text = re.sub(r"http\S+", "", text) # replaces URLs starting with http
5     text = re.sub(r"www.\S+", "", text) # replaces URLs starting with wwe
6     return text
```

```
1 text = 'https://thistutorialisawesome.com is best site'
2 print(f'Before - {text}')
3 print(f'After  - {removeURLs(text)}')
```

```
Before - https://thistutorialisawesome.com is best site
After  -  is best site
```

## 3) Removal of punctuations

**Punctuation don't add any value to overall text and can be removed without impacting**

```
1 import string
2 string.punctuation # checking punctuations
```

```
'!"#$%&\'()*+,-./:;<=>?@[\\]^_`{|}~'
```

```
1 def removePunctuation(text):
2     return "".join([char for char in text if char not in string.punctuation])
```

```
1 text = 'Is that seriously how you spell my name? I am not joking.'
2 print(f'Before - {text}')
3 print(f'After  - {removePunctuation(text)}')
```

```
Before - Is that seriously how you spell my name? I am not joking.
After  - Is that seriously how you spell my name I am not joking
```

## 4) Removal of stopwords

**In sentiment analysis, stopwords like pronouns can overweight actual words expressing emotions, which leads to poor for perfomance of model. So stopwords should be removed as a part of text preprocessing**

```
1 import nltk
2 nltk.download('stopwords')
3 from nltk.corpus import stopwords
4 stopword = nltk.corpus.stopwords.words('english')
5 print(stopword[0:500:25])
```

```
['i', 'herself', 'been', 'with', 'here', 'very', 'doesn', 'won']
[nltk_data] Downloading package stopwords to
[nltk_data]     /Users/afsarequebal/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
```

```
1 def removeStopwords(text):
2     return " ".join([word for word in re.split('\W+', text)
3         if word not in stopword])
```

```
1 text = 'Is that seriously how you spell my name? I am not joking.'
2 print(f'Before - {text.lower()}')
3 print(f'After  - {removeStopwords(text.lower())}')
```

```
Before - is that seriously how you spell my name? i am not joking.
After  - seriously spell name joking
```

## 5) Tokenization - It is essentially splitting a phrase, sentence, paragraph, or an entire text document into smaller units, such as individual words or terms. Each of these smaller units are called tokens.

**NLTK contains a module called tokenize() which further classifies into two sub-categories: Word tokenize: We use the word_tokenize() method to split a sentence into tokens or words. Sentence tokenize: We use the sent_tokenize() method to split a document or paragraph into sentences.**

```
1 nltk.download('punkt')
2 from nltk.tokenize import sent_tokenize, word_tokenize
3 def sentenceTokenize(text):
4     return sent_tokenize(text)
5
```

```
6 def wordTokenize(text):
7     return word tokenize(text)
```

```
[nltk_data] Downloading package punkt to
[nltk_data]     /Users/afsarequebal/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
```

```
1 text = 'Its a wonderful day in Boston. I really like NLP. I will be visting some nearby park today.'
2 print(f'Before - {text}')
3 print(f'After  - {sentenceTokenize(text.lower())}')
```

```
Before - Its a wonderful day in Boston. I really like NLP. I will be visting some nearby
After  - ['its a wonderful day in boston.', 'i really like nlp.', 'i will be visting som
```

```
1 text = 'Its a wonderful day in Boston. I really NLP. I will be visting some nearby park today.'
2 print(f'Before - {text}')
3 print(f'After  - {wordTokenize(text.lower())}')
```

```
Before - Its a wonderful day in Boston. I really NLP. I will be visting some nearby park
After  - ['its', 'a', 'wonderful', 'day', 'in', 'boston', '.', 'i', 'really', 'nlp', '.'
```

## ⌄ 6) Stemming

**Stemming is the process of reducing words to their stem/roots. It help in removing redundant words. For eg. root of connect, connected, connecting, connection is same. We will be using NLTK to perform this task**

```
1 nltk.download('wordnet')
2 import nltk
3 from nltk.stem import PorterStemmer
4 from nltk.tokenize import word_tokenize
5 wn = nltk.WordNetLemmatizer()
6 ps = nltk.PorterStemmer()
```

```
[nltk_data] Downloading package wordnet to
[nltk_data]     /Users/afsarequebal/nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
```

```
1 def performStemming(text):
2     return" ".join([ps.stem(word) for word in re.split('\W+', text)])
```

```
1 text = 'It am troubled by riding bicyle daily. I have no troubles taking a ride in lyft to college.'
2 print(f'Before - {text}')
3 print(f'After  - {performStemming(text.lower())}')
```

```
Before - It am troubled by riding bicyle daily. I have no troubles taking a ride in lyft
After  - it am troubl by ride bicyl daili i have no troubl take a ride in lyft to colleg
```

## ∨  7) Lemmatization

**It is similar to stemming except the lemmatized word belongs to the language. It also allows us to specify verb or noun to be used as parameter.**

```
1 def performLemmatization(text):
2     return" ".join([wn.lemmatize(word,'v') for word in re.split('\W+', text)])
```

```
1 text = 'It am troubled by riding bicyle daily. I have no troubles taking a ride in lyft to college.'
2 print(f'Before - {text}')
3 print(f'After  - {performLemmatization(text.lower())}')
```

```
Before - It am troubled by riding bicyle daily. I have no troubles taking a ride in lyft
After  - it be trouble by rid bicyle daily i have no trouble take a ride in lyft to coll
```

Exercise:

Write a Python program using NLTK to:

1. Tokenize the following text into words:

I am a fourth-year bachelor's student in Data Science and AI, an exciting field that combines cutting-edge technology with the power of data to solve real-world problems!

2. Remove stopwords from the tokenized words.
3. Apply **stemming** to the filtered words using the Porter Stemmer.
4. Apply **lemmatization** to the filtered words..
5. Explain the difference between stemming and lemmatization using your results.