

Assignment 7

Tokenization

```
In [1]: 1 import nltk
        2 nltk.download('punkt')
        3 nltk.download('wordnet')
        4 nltk.download('averaged_perceptron_tagger')
        5 nltk.download('stopwords')
        6 from nltk import sent_tokenize
        7 from nltk import word_tokenize
        8 from nltk.corpus import stopwords
        9
       10
```

```
[nltk_data] Downloading package punkt to
[nltk_data]   C:\Users\omraj\AppData\Roaming\nltk_data...
[nltk_data]   Unzipping tokenizers\punkt.zip.
[nltk_data] Downloading package wordnet to
[nltk_data]   C:\Users\omraj\AppData\Roaming\nltk_data...
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]   C:\Users\omraj\AppData\Roaming\nltk_data...
[nltk_data]   Unzipping taggers\averaged_perceptron_tagger.zip.
[nltk_data] Downloading package stopwords to
[nltk_data]   C:\Users\omraj\AppData\Roaming\nltk_data...
[nltk_data]   Unzipping corpora\stopwords.zip.
```

```
In [2]: 1 text='Real madrid is set to win the UCL for the season . Benzema might win B
```

```
In [3]: 1 tokens_sents = nltk.sent_tokenize(text)
        2 print(tokens_sents)
```

```
['Real madrid is set to win the UCL for the season .', 'Benzema might win Balon
dor .', 'Salah might be the runner up']
```

```
In [4]: 1 tokens_words = nltk.word_tokenize(text)
        2 print(tokens_words)
```

```
['Real', 'madrid', 'is', 'set', 'to', 'win', 'the', 'UCL', 'for', 'the', 'seaso
n', '.', 'Benzema', 'might', 'win', 'Balon', 'dor', '.', 'Salah', 'might', 'b
e', 'the', 'runner', 'up']
```

```
In [5]: 1 from nltk.stem import PorterStemmer
        2 from nltk.stem.snowball import SnowballStemmer
        3 from nltk.stem import LancasterStemmer
```

```
In [6]: 1 stem=[]
        2 for i in tokens_words:
        3     ps = PorterStemmer()
        4     stem_word= ps.stem(i)
        5     stem.append(stem_word)
        6 print(stem)
        7
        8
```

```
['real', 'madrid', 'is', 'set', 'to', 'win', 'the', 'ucl', 'for', 'the', 'seaso
n', '.', 'benzema', 'might', 'win', 'balon', 'dor', '.', 'salah', 'might', 'b
e', 'the', 'runner', 'up']
```

Lemmatization

```
In [7]: 1 import nltk
        2 from nltk.stem import WordNetLemmatizer
        3 lemmatizer = WordNetLemmatizer()
```

```
In [8]: 1 lemmatized_output = ' '.join([lemmatizer.lemmatize(w) for w in stem])
        2 print(lemmatized_output)
```

```
real madrid is set to win the ucl for the season . benzema might win balon dor
. salah might be the runner up
```

```
In [9]: 1 leme=[]
        2 for i in stem:
        3     lemetized_word=lemmatizer.lemmatize(i)
        4     leme.append(lemetized_word)
        5 print(leme)
```

```
['real', 'madrid', 'is', 'set', 'to', 'win', 'the', 'ucl', 'for', 'the', 'seaso
n', '.', 'benzema', 'might', 'win', 'balon', 'dor', '.', 'salah', 'might', 'b
e', 'the', 'runner', 'up']
```

Part of Speech Tagging

```
In [10]: 1 print("Parts of Speech: ",nltk.pos_tag(leme))
        2
```

```
Parts of Speech: [('real', 'JJ'), ('madrid', 'NN'), ('is', 'VBZ'), ('set', 'VB
N'), ('to', 'TO'), ('win', 'VB'), ('the', 'DT'), ('ucl', 'NN'), ('for', 'IN'),
('the', 'DT'), ('season', 'NN'), ('.', '.'), ('benzema', 'NN'), ('might', 'M
D'), ('win', 'VB'), ('balon', 'NN'), ('dor', 'NN'), ('.', '.'), ('salah', 'N
N'), ('might', 'MD'), ('be', 'VB'), ('the', 'DT'), ('runner', 'NN'), ('up', 'R
P')]
```

Stop Word

```
In [11]: 1 sw_nltk = stopwords.words('english')
          2 print(sw_nltk)
```

```
['a', 'about', 'above', 'after', 'again', 'against', 'ain', 'all', 'am', 'an',
'and', 'any', 'are', 'aren', "aren't", 'as', 'at', 'be', 'because', 'been', 'be
fore', 'being', 'below', 'between', 'both', 'but', 'by', 'can', 'couldn', "coul
dn't", 'd', 'did', 'didn', "didn't", 'do', 'does', 'doesn', "doesn't", 'doing',
'don', "don't", 'down', 'during', 'each', 'few', 'for', 'from', 'further', 'ha
d', 'hadn', "hadn't", 'has', 'hasn', "hasn't", 'have', 'haven', "haven't", 'hav
ing', 'he', "he'd", "he'll", 'her', 'here', 'hers', 'herself', "he's", 'him',
'himself', 'his', 'how', 'i', "i'd", 'if', "i'll", "i'm", 'in', 'into', 'is',
'isn', "isn't", 'it', "it'd", "it'll", "it's", 'its', 'itself', "i've", 'just',
'll', 'm', 'ma', 'me', 'mightn', "mightn't", 'more', 'most', 'mustn', "must
n't", 'my', 'myself', 'needn', "needn't", 'no', 'nor', 'not', 'now', 'o', 'of',
'off', 'on', 'once', 'only', 'or', 'other', 'our', 'ours', 'ourselves', 'out',
'over', 'own', 're', 's', 'same', 'shan', "shan't", 'she', "she'd", "she'll",
"she's", 'should', 'shouldn', "shouldn't", "should've", 'so', 'some', 'such',
't', 'than', 'that', "that'll", 'the', 'their', 'theirs', 'them', 'themselves',
'then', 'there', 'these', 'they', "they'd", "they'll", "they're", "they've", 't
his', 'those', 'through', 'to', 'too', 'under', 'until', 'up', 've', 'very', 'w
as', 'wasn', "wasn't", 'we', "we'd", "we'll", "we're", 'were', 'weren', "were
n't", "we've", 'what', 'when', 'where', 'which', 'while', 'who', 'whom', 'why',
'will', 'with', 'won', "won't", 'wouldn', "wouldn't", 'y', 'you', "you'd", "yo
u'll", 'your', "you're", 'yours', 'yourself', 'yourselves', "you've"]
```

```
In [12]: 1 words = [word for word in text.split() if word.lower() not in sw_nltk]
          2 new_text = " ".join(words)
          3 print(new_text)
```

Real madrid set win UCL season . Benzema might win Balon dor . Salah might runn
er