

In [1]:

```
pip install nltk
```

Requirement already satisfied: nltk in c:\users\ahmad\anaconda3\lib\site-packages (3.4.5)

Requirement already satisfied: six in c:\users\ahmad\anaconda3\lib\site-packages (from nltk) (1.12.0)

Note: you may need to restart the kernel to use updated packages.

In [2]:

```
import nltk
```

In [3]:

```
nltk.download('punkt')
```

[nltk_data] Downloading package punkt to

[nltk_data] C:\Users\ahmad\AppData\Roaming\nltk_data...

[nltk_data] Package punkt is already up-to-date!

Out[3]:

True

In [4]:

```
text=" Welcome readers. I hope you find it interesting. Please do reply."
```

In [5]:

```
from nltk.tokenize import sent_tokenize
```

In [6]:

```
sent_tokenize
```

Out[6]:

```
<function nltk.tokenize.sent_tokenize(text, language='english')>
```

In [7]:

```
sent_tokenize(text)
```

Out[7]:

```
[' Welcome readers.', 'I hope you find it interesting.', 'Please do reply.']
```

In [8]:

```
tokenizer = nltk.data.load("tokenizers/punkt/english.pickle")
```

In [9]:

```
text=" Welcome readers. I hope you find it interesting. Please do reply."  
tokenizer.tokenize(text)
```

Out[9]:

```
[' Welcome readers.', 'I hope you find it interesting.', 'Please do repl  
y.']
```

In []:

In [10]:

```
Arabic_text="مرحبا بكم. نحن نتعلم اساسيات مبادئ استرجاع المعلومات."  
tokenizer.tokenize(Arabic_text)
```

Out[10]:

```
['مرحبا بكم.', 'نحن نتعلم اساسيات مبادئ استرجاع المعلومات']
```

In [11]:

```
text=nltk.word_tokenize("Welcome readers. I hope you find it interesting. Please do rep  
ly..»")
```

In [12]:

```
print (text)
```

```
['Welcome', 'readers', '.', 'I', 'hope', 'you', 'find', 'it', 'interestin  
g', '.', 'Please', 'do', 'reply..', '»']
```

In [13]:

```
input(text)
```

Out[13]:

```
'hi'
```

In [14]:

```
import nltk  
from nltk.tokenize import TreebankWordTokenizer
```

In [15]:

```
tokenizer = TreebankWordTokenizer()  
tokenizer.tokenize("Have a nice day. You do great!")
```

Out[15]:

```
['Have', 'a', 'nice', 'day.', 'You', 'do', 'great', '!']
```

In [16]:

```
Arabic=input("Please write a text")
```

In [17]:

```
Arabic_tokenizer=nltk.word_tokenize(Arabic)

print (Arabic_tokenizer)

['I', 'am', 'Ahmad']
```

In [18]:

```
from nltk.tokenize import RegexpTokenizer
```

In [19]:

```
tokenizer=RegexpTokenizer("[\w]+")
```

In [20]:

```
tokenizer.tokenize("Don't hesitate to ask questions or send to me your question to mohs  
arem@gmail.com")
```

Out[20]:

```
['Don',  
't',  
'hesitate',  
'to',  
'ask',  
'questions',  
'or',  
'send',  
'to',  
'me',  
'your',  
'question',  
'to',  
'mohsarem',  
'gmail',  
'com']
```

In [21]:

```
tokenizer=RegexpTokenizer("\S+@\S+")
```

In [22]:

```
tokenizer.tokenize("Don't hesitate to ask questions or send to me your question to mohs  
arem@gmail.com")
```

Out[22]:

```
['mohsarem@gmail.com']
```

In [23]:

```
text=[" It is a pleasant evening.", "Guests, who came from US arrived at the venue", "Food  
was tasty."]
```

In [24]:

```
from nltk.tokenize import word_tokenize
tokenized_docs=[word_tokenize(doc) for doc in text]
print(tokenized_docs)
```

```
[['It', 'is', 'a', 'pleasant', 'evening', '.'], ['Guests', ',', 'who', 'came', 'from', 'US', 'arrived', 'at', 'the', 'venue'], ['Food', 'was', 'tasty', '.']]
```

In [25]:

```
import re
import string
x=re.compile('[%s]' % re.escape(string.punctuation))
```

In [28]:

```
tokenized_docs_no_punctuation = []
for review in tokenized_docs:
    new_review = []
    for token in review:
        new_token = x.sub(u'', token)
        if not new_token == u'':
            new_review.append(new_token)
    tokenized_docs_no_punctuation.append(new_review)
print(tokenized_docs_no_punctuation)
```

```
[['It', 'is', 'a', 'pleasant', 'evening'], ['Guests', 'who', 'came', 'from', 'US', 'arrived', 'at', 'the', 'venue'], ['Food', 'was', 'tasty']]
```

In [29]:

```
print(text[0].upper())
```

IT IS A PLEASANT EVENING.

In [30]:

```
print(text[0].lower())
```

it is a pleasant evening.

In [31]:

```
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords
stops=set(stopwords.words('english'))
words=["Don't", 'hesitate', 'to', 'ask', 'questions']
[word for word in words if word not in stops]
```

```
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\ahmad\AppData\Roaming\nltk_data...
[nltk_data] Unzipping corpora\stopwords.zip.
```

Out[31]:

```
["Don't", 'hesitate', 'ask', 'questions']
```

```
Text= "NLTK allows you to convert Text into Lowercase and uppercase. Don't hesitate to askquestions"
```

```
print(stopwords.words('english'))
```

```
print(stopwords.words('Arabic'))
```

5/6

In []: