-	['هل', 'القام', 'المواجهات', 'المباشرة', 'لمباشرة', 'لمباشرة', 'المباشرة', 'لمباشرة', 'لمبا
=	Text Normalizaion (Cleaning)
	How to remove all this from the text?
	 Stop Words: (such as "the", "a", "an", "in") — a search engine has been programmed to ignore. Punctuation: [", "", #", \$5, "%, "&", "", ", ", ", ", ", ", ", ", ", ", ",
	• The key concept here is that stemming sometime destroy the word unlike lemmatization where we keep the meaning.
	stop words in EN
	#NLTK(Natural Language Toolkit) in python has a list of stopwords stored in 16 different languages from nltk.corpus import stopwords from nltk.tokenize import word_tokenize
l u s t e	print(stopwords.words('english')) ('a', 'about', 'above', 'after', 'again', 'against', 'ain', 'all', 'am', 'an', 'and', 'any', 'are', 'aren', "aren't", 'as', 'at', 'be', 'because', 'been', 'before', 'being', 'below', 'between', 'both', 'but', 'by', 'can', 'couldn', 'didn', 'didn', "didn't", 'do', 'does', 'doesn', "doesn't", 'doing', 'don', "don't", 'down', 'during', 'each', 'few', 'for', 'from', 'further', 'had', 'hadn', "hadn't", 'has', 'hasn', "hasn't", 'have', 'haven', having', 'he', "he'd", "he'll", 'her', 'here', 'hers', 'herself', "he's", 'him', 'himself', 'his', 'how', 'i', "i'll", "i'll", "i'm", 'into', 'is', 'isn', "isn't", 'it', "it'd", "it'll", "it's", 'itself', "ive", 'it', 'll', "m', 'ma, 'me', 'mightn't", 'more', 'most', 'mustn't", 'my', 'myself', 'needn't", 'no', 'nor', 'not', 'now', 'o', 'of', 'off', 'on', 'onee', 'only', 'or', 'other', 'our', 'ourselves', 'o', 'over', 'own', 're', 's', 'same', 'shan', "shan't", 'she', "she'd", "she's", 'should', 'shouldn't", "should've", 'so', 'some', 'such', 't', 'than', 'that', "that'll", 'the', 'their', 'theirs', 'them', 'them', 'ss', 'there', 'there', 'these', 'they', "they'd", "they'd", "they're", "they've", 'those', 'through', 'to', 'too', 'under', 'until', 'up', 've', 'very', 'wasn', "wasn't", 'we', "yourself', 'yourselves', "you've" 'weren't", "we've", 'what', 'where', 'which', 'which', 'whom', 'why', 'will', 'won', "won't", 'wouldn't", "you'd", "you'd", "you'll", 'your', "you're", 'yourself', 'yourselves', "you've"
	Task 1 : Removing stop words with NLTK ?
[]:	
[]:	
	<pre>example_sent = """This is a sample sentence, showing off the stop words filtration.""" stop_words = stopwords.words('english') word_tokens = word_tokenize(example_sent) filtered_sentence = [] for w in word_tokens: if w.lower() not in stop_words: filtered_sentence.append(w.lower())</pre>
[['This', 'is', 'a', 'sample', 'sentence', ',', 'showing', 'off', 'the', 'stop', 'words', 'filtration', '.'] ['sample', 'sentence', ',', 'showing', 'stop', 'words', 'filtration', '.'] [sample sentence , showing stop words filtration .
	stop words in AR
	<pre>stop_word_arabic = set(stopwords.words("arabic")) print(len(stop_word_arabic))</pre>
7	
	"
[17]:	from string import punctuation
	<pre>punctuation = list(punctuation) print(punctuation)</pre>
	['!', '"', '\$', '\$', '\$', '&', "'", '(', ')', '*', '+', ',', '-', '.', '/', ';', '<', '=', '>', '@', '[', '\\', ']', '^', '\\','\\', '\\','\\', '\\','\\','\\','\\','\\','\\','\\','\\','\\','\\
	Task 2 : Removing punctuation with tokenization ?
[]:	
[]:	
[]:	
	<pre>example_sent = """My Email address is: taneshbalodi8@gmail.com.""" word_tokens = word_tokenize(example_sent) filtered_sentence = []</pre>
	<pre>for w in word_tokens: if w not in punctuation: filtered_sentence.append(w)</pre>
	<pre>print (word_tokens) print (filtered_sentence)</pre>
[['My', 'Email', 'address', 'is', ':', 'taneshbalodi8', '@', 'gmail.com', '.'] ['My', 'Email', 'address', 'is', 'taneshbalodi8', 'gmail.com']
L	Task 3 : Removing punctuation without tokenization
[]:	racit o . Itomoving pariotation without tokemzation
[]:	

Text Normalization

['Hello', 'Mr.', 'Smith', ',', 'how', 'are', 'you', 'doing', 'today', '?', 'The', 'awesome', '.', 'The', 'sky', 'is', 'pinkish-blue', '.', 'You', 'should', "n't", 'eat', 'cardboar

It involves cleaning and preprocessing text data to make it consistent and usable for different NLP tasks. The process includes a variety of techniques, such as case normalization, punctuation removal, stop word removal, stemming, and lemmatization.

Tokenization

1.1. Tokenization Using NLTK EN1.2. Tokenization Using NLTK AR

Text Normalization (Cleaning)
 2.1. Stop words in EN

2.3. Stemming using nltk2.4. Stemming in AR

2.5. Lemmatization using nltk2.6. Lemmatization in AR

Import Library (nltk)

Ref: "https://www.nltk.org/data.html"

In [39]: # to know the functions in the library

'AbstractLazySequence',

'ApplicationExpression',

'BigramAssocMeasures']

Tokenization

Tokenization Using NLTK EN

In [69]: from nltk.tokenize import word_tokenize
#word_tokenize is a function

print (word_tokenize (EXAMPLE_TEXT))

Difference between split and token

print('----')

print (word_tokenize(text))

In [71]: from nltk.tokenize import sent_tokenize

Tokenization Using NLTK AR

print (word_tokenize(EXAMPLE_TEXT))

print('----')
print(word_tokenize(EXAMPLE_TEXT))

print(sent_tokenize(s))
#whitespaceTokenizer?

In [74]: print(EXAMPLE_TEXT.split())

EXAMPLE_TEXT = """

print(text.split())

d', '.']

dir(nltk)[0:10]

'ARLSTem2',

'AffixTagger',
'AlignedSent',
'Alignment',
'AnnotationTask',

'Assignment',

nltk.download() # to download all features

In [68]: import nltk

Out[39]: ['ARLSTem',

Out[68]: True

• NLTK Corpora (DataSets) : "https://www.nltk.org/nltk_data/"

■ 2.7. Special arabic cleaning functions

• work with human language data for applying in statistical natural language processing (NLP).

showing info https://raw.githubusercontent.com/nltk/nltk data/gh-pages/index.xml

Tokenization: mean split the sentence into tokens(words) and called Word tokenizer

• Apply tokenization Text -> This product is amazing, but the delivery was late.

Hello Mr. Smith, how are you doing today? The weather is great,

print ('======="")

print ('======"")

In [19]: example_sent = """My Email address is: hossamfares180100@gmail.com."""

words = ["program", "programs", "programmer", "programming", "programmers"]

filtered_sentence = []

for w in example_sent:

print (example_sent)
print (filtered_sentence)

if w not in punctuation:

print(''.join(filtered_sentence))

stemming using nltk

In [21]: **from** nltk.stem **import** PorterStemmer

ps = PorterStemmer()

for w in words:

filtered_sentence.append(w)

My Email address is: hossamfares180100@gmail.com.

My Email address is hossamfares180100gmailcom

from nltk.tokenize import word_tokenize

choose some words to be stemmed

print(w, " : ", ps.stem(w))

print(lemmatizer.lemmatize(word, 'v'))

lemmatization in AR

for word in words :

In [32]: words = ['ایجری','یجرون','بیرون','بریی']

print(lemmatizer.lemmatize(word))

Task 4: Create Method Take (Text as input) then: (15 MIN)

if word.lower() not in stop_words and w not in punctuation :

lemmatizer_text = [lemmatizer.lemmatize(word) for word in filtered_sentence]

clean_text("Welcome to the NLP Course!. Natural Language Processing (NLP) is an exciting field of Artificial Intelligence that enables computers to understand, interpret, and generate human language.")

filtered_sentence.append(word.lower())

ps_text = [ps.stem(word) for word in filtered_sentence]

returns a cleaned version by removing stopwords, punctuation, applying lemmatization, and stemming:

be be be be

ا لجری تجری

جر ي يـجـر ي

In [114... **def** clean_text(text):

Out[118... (['welcom',

'nlp',
'course!.',
'natur',
'languag',
'process',
'(nlp)',
'excit',
'field',
'artifici',
'intellig',
'enabl',
'comput',
'understand,',
'interpret,',

'gener', 'human', 'language.'], ['welcome', 'nlp', 'course!.', 'natural', 'language', 'processing', '(nlp)', 'exciting', 'field', 'artificial', 'intelligence', 'enables', 'computer', 'understand,', 'interpret,', 'generate', 'human', 'language.'])

Extra Topics:

In [33]: pip install qalsadi

Collecting qalsadi

Some works in Arabic

Downloading qalsadi-0.5-py3-none-any.whl.metadata (12 kB)

Downloading alyahmor-0.2-py3-none-any.whl.metadata (11 kB)

Preparing metadata (setup.py): finished with status 'done'

Downloading Naftawayh-0.4-py3-none-any.whl.metadata (8.7 kB)

Preparing metadata (setup.py): finished with status 'done'

Downloading PyArabic-0.6.15-py3-none-any.whl.metadata (10 kB)

Downloading Tashaphyne-0.3.6-py3-none-any.whl.metadata (18 kB)

Downloading aiofiles-24.1.0-py3-none-any.whl.metadata (10 kB)

----- 0.0/264.3 kB ? eta -:--:-

----- 0.0/65.1 kB ? eta -:--:-

----- 0.0/360.5 kB ? eta -:--:-

----- 0.0/6.3 MB ? eta -:--:-

----- 65.1/65.1 kB 1.2 MB/s eta 0:00:00

----- 0.1/6.3 MB 2.3 MB/s eta 0:00:03 - ----- 0.2/6.3 MB 2.1 MB/s eta 0:00:03 - ----- 0.3/6.3 MB 2.2 MB/s eta 0:00:03 -- ----- 0.4/6.3 MB 2.2 MB/s eta 0:00:03 --- 0.5/6.3 MB 2.4 MB/s eta 0:00:03 --- 0.6/6.3 MB 2.2 MB/s eta 0:00:03 ---- 0.7/6.3 MB 2.3 MB/s eta 0:00:03 ---- 0.8/6.3 MB 2.3 MB/s eta 0:00:03 ---- 0.9/6.3 MB 2.3 MB/s eta 0:00:03 ----- 1.0/6.3 MB 2.3 MB/s eta 0:00:03 ----- 1.1/6.3 MB 2.3 MB/s eta 0:00:03 ----- 1.2/6.3 MB 2.3 MB/s eta 0:00:03 ----- 1.3/6.3 MB 2.3 MB/s eta 0:00:03 ----- 1.4/6.3 MB 2.2 MB/s eta 0:00:03 ----- 1.5/6.3 MB 2.3 MB/s eta 0:00:03 ----- 1.6/6.3 MB 2.3 MB/s eta 0:00:03 ------ 1.8/6.3 MB 2.3 MB/s eta 0:00:03 ----- 1.9/6.3 MB 2.3 MB/s eta 0:00:02 ----- 2.0/6.3 MB 2.3 MB/s eta 0:00:02 ----- 2.1/6.3 MB 2.3 MB/s eta 0:00:02 ----- 2.2/6.3 MB 2.3 MB/s eta 0:00:02 ----- 2.3/6.3 MB 2.3 MB/s eta 0:00:02 ----- 2.4/6.3 MB 2.3 MB/s eta 0:00:02 ----- 2.5/6.3 MB 2.3 MB/s eta 0:00:02 ----- 2.6/6.3 MB 2.3 MB/s eta 0:00:02 ----- 2.7/6.3 MB 2.3 MB/s eta 0:00:02 ----- 2.8/6.3 MB 2.3 MB/s eta 0:00:02 ----- 2.9/6.3 MB 2.3 MB/s eta 0:00:02 ----- 3.0/6.3 MB 2.3 MB/s eta 0:00:02 ----- 3.1/6.3 MB 2.3 MB/s eta 0:00:02 ----- 3.2/6.3 MB 2.3 MB/s eta 0:00:02 ----- 3.3/6.3 MB 2.3 MB/s eta 0:00:02 ----- 3.4/6.3 MB 2.3 MB/s eta 0:00:02 ----- 3.5/6.3 MB 2.3 MB/s eta 0:00:02 ----- 3.6/6.3 MB 2.3 MB/s eta 0:00:02 ----- 3.8/6.3 MB 2.3 MB/s eta 0:00:02 ----- 3.9/6.3 MB 2.3 MB/s eta 0:00:02 ----- 4.0/6.3 MB 2.3 MB/s eta 0:00:02 ----- 4.0/6.3 MB 2.3 MB/s eta 0:00:02 ------ 4.0/6.3 MB 2.3 MB/s eta 0:00:02 ----- 4.1/6.3 MB 2.0 MB/s eta 0:00:02 ----- 4.8/6.3 MB 2.3 MB/s eta 0:00:01 ----- 4.9/6.3 MB 2.3 MB/s eta 0:00:01 ----- 5.0/6.3 MB 2.3 MB/s eta 0:00:01 ----- 5.1/6.3 MB 2.3 MB/s eta 0:00:01 ----- 5.2/6.3 MB 2.3 MB/s eta 0:00:01 ----- 5.3/6.3 MB 2.3 MB/s eta 0:00:01 ----- 5.4/6.3 MB 2.3 MB/s eta 0:00:01 ----- 5.6/6.3 MB 2.3 MB/s eta 0:00:01 ----- 5.6/6.3 MB 2.3 MB/s eta 0:00:01 ----- 5.8/6.3 MB 2.3 MB/s eta 0:00:01 ----- 5.9/6.3 MB 2.3 MB/s eta 0:00:01 ----- -- 6.0/6.3 MB 2.3 MB/s eta 0:00:01 ----- - 6.1/6.3 MB 2.3 MB/s eta 0:00:01 ----- - 6.2/6.3 MB 2.3 MB/s eta 0:00:01 ----- 6.3/6.3 MB 2.3 MB/s eta 0:00:01 ----- 6.3/6.3 MB 2.3 MB/s eta 0:00:01 ----- 6.3/6.3 MB 2.2 MB/s eta 0:00:00

Downloading libqutrub-1.2.4.1-py3-none-any.whl.metadata (7.5 kB)

Downloading mysam_tagmanager-0.4-py3-none-any.whl.metadata (10 kB)

Downloading Arabic_Stopwords-0.4.3-py3-none-any.whl.metadata (8.9 kB)

Downloading arramooz_pysqlite-0.4.2-py3-none-any.whl.metadata (4.0 kB)

----- 0.0/46.1 kB ? eta -:--:-

----- 30.7/46.1 kB 1.3 MB/s eta 0:00:01 ----- 46.1/46.1 kB 568.6 kB/s eta 0:00:00

Requirement already satisfied: or json in c:\users\fady\anaconda3\lib\site-packages (from pickledb>=0.9.2->qalsadi) (3.10.5)

Requirement already satisfied: six>=1.14.0 in c:\users\fady\anaconda3\lib\site-packages (from pyarabic>=0.6.7->qalsadi) (1.16.0)

Collecting Arabic-Stopwords>=0.4.2 (from qalsadi)

Collecting arramooz-pysqlite>=0.4.2 (from qalsadi)

Downloading codernitydb3-0.6.0.tar.gz (46 kB)

Collecting mysam-tagmanager>=0.3.3 (from galsadi)

Collecting alyahmor>=0.2 (from galsadi)

Collecting codernitydb3 (from qalsadi)

Preparing metadata (setup.py): started

Collecting libqutrub>=1.2.3 (from qalsadi)

Collecting naftawayh>=0.3 (from qalsadi)

Collecting pickledb>=0.9.2 (from qalsadi)
Downloading pickledb-1.3.2.tar.gz (10 kB)
Preparing metadata (setup.py): started

Collecting pyarabic>=0.6.7 (from qalsadi)

Collecting tashaphyne>=0.3.4.1 (from qalsadi)

Collecting aiofiles (from pickledb>=0.9.2->galsadi)

Downloading qalsadi-0.5-py3-none-any.whl (264 kB)

Downloading alyahmor-0.2-py3-none-any.whl (65 kB)

Downloading Arabic_Stopwords-0.4.3-py3-none-any.whl (360 kB)

Downloading arramooz_pysqlite-0.4.2-py3-none-any.whl (6.3 MB)

Downloading libqutrub-1.2.4.1-py3-none-any.whl (138 kB)

Downloading mysam_tagmanager-0.4-py3-none-any.whl (37 kB)

Downloading Naftawayh-0.4-py3-none-any.whl (332 kB)

Downloading PyArabic-0.6.15-py3-none-any.whl (126 kB)

Downloading Tashaphyne-0.3.6-py3-none-any.whl (251 kB)

Downloading aiofiles-24.1.0-py3-none-any.whl (15 kB)

Building wheel for pickledb (setup.py): started

Successfully built pickledb codernitydb3

lemmer = qalsadi.lemmatizer.Lemmatizer()
words = ['یجری','یجری','یبرون','یبری',']

print (st.stem('یقول'),end=' ')
print (st.stem('یقولون'),end=' ')
print (st.stem('تقول'),end=' ')

lemmas = lemmer.lemmatize_text(text)

In [36]: lemmas = lemmer.lemmatize_text(text, return_pos=True)

print (st.stem('مقوله'))

جرة تجرة جرى جرة جرى يقل يقل تقل قول

print(lemmas)

print(lemmas)

In [37]: # put diacritics on text

print(lemmas)

In []: #

lemmer.set_vocalized_lemma()

In [38]: # special cleaning to arabic text

def normalize_arabic(text):

def remove_diacritics(text):

def remove_repeating_char(text):

return text

text = re.sub("[||||||", "|", text)

text = re.sub("g", "g", text)

text = re.sub("s", "g", text)

text = re.sub("s", "g", text)

text = re.sub("o", "g", text)

text = re.sub("a", "g", text)

text = re.sub("a", "g", text)

text = re.sub(arabic_diacritics, '', text)

return re.sub(r'(.)\1+', r'\1', text)

print(remove_diacritics(' '.join(lemmas)))

print(normalize_arabic(remove_diacritics(' '.join(lemmas))))

arabic_diacritics = re.compile("""

lemmas = lemmer.lemmatize_text(text)

special arabic cleaning functions

| # Tashdid | # Fatha

[| # Damma

| # Kasra

""", re.VERBOSE)

| # Tanwin Fath

| # Tanwin Damm

| # Tanwin Kasr | # Sukun | # Tatwil/Kashida

print(lemmer.lemmatize(word),end=' ')

In [34]: **import** qalsadi.lemmatizer

for word in words :

print()

Building wheel for codernitydb3 (setup.py): started

Building wheels for collected packages: pickledb, codernitydb3

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

[notice] A new release of pip is available: 24.1 -> 25.0.1
[notice] To update, run: python.exe -m pip install --upgrade pip

----- 0.0/139.0 kB ? eta -:--:-

----- 0.0/332.6 kB ? eta -:--:-

----- 0.0/126.4 kB ? eta -:--:-

----- 0.0/251.5 kB ? eta -:--:--

Building wheel for pickledb (setup.py): finished with status 'done'

Building wheel for codernitydb3 (setup.py): finished with status 'done'

------ 92.2/139.0 kB 2.6 MB/s eta 0:00:01 ------ 139.0/139.0 kB 2.0 MB/s eta 0:00:00

----- 92.2/126.4 kB 2.6 MB/s eta 0:00:01 ----- 126.4/126.4 kB 2.5 MB/s eta 0:00:00

Created wheel for pickledb: filename=pickleDB-1.3.2-py3-none-any.whl size=5307 sha256=104132bbb93fc846574c9be4c58e42d679a1fbbf1a9850ea2a9bd3a3f1bff7e3

Created wheel for codernitydb3: filename=codernitydb3-0.6.0-py3-none-any.whl size=59877 sha256=4ce7e2e44fbe5f3a8d65b74a8392b4ffc456975e4be3e6a6cc7e4bbeaacdae8e

Installing collected packages: mysam-tagmanager, pyarabic, codernitydb3, aiofiles, tashaphyne, pickledb, libqutrub, arramooz-pysqlite, Arabic-Stopwords, naftawayh, alyahmor, qalsadi

Successfully installed Arabic-Stopwords-0.4.3 aiofiles-24.1.0 alyahmor-0.2 arramooz-pysqlite-0.4.2 codernitydb3-0.6.0 libqutrub-1.2.4.1 mysam-tagmanager-0.4 naftawayh-0.4 pickledb-1.3.2 pyarabic-0.6.15 qalsadi-0.5 tashaphyne-0.3.6

ب الملك؟ اللغة "الكلاسيكية" (الفصحى) موجودة في كل اللغات وكذلك اللغة "الدارجة" .. الفرنسية التين ندرس في المدرسة ليست الفرنسية التين يستخدمها الناس في شوارع باريس .. وملكة بريطانيا لا تخطب بلغة شوارع لندن .. لكل مقام مقال"" = 153]

'noun'), ('امقام', 'noun'), ('الا'), 'stopword'), ('الاارع', 'stopword'), ('الاارع', 'noun'), ('الدان', 'noun'), ('الدان', 'verb'), ('...', 'pounct'), ('الله'), 'stopword'), ('الله ', 'noun'), (''لله ',

مَلْ إِحْتَاجَ إِلَى تَرْجَمَةً كَيْ تَفَهُمْ خَطَّابٌ مَلَكْ ؟ لُغَةٌ " كِلاَسِيكِيّْ " (فُضْحَى) مَوْجُودْ فِي كُلِّ لُغَةٌ " دَارِجٌ ".. فَرَنْسِيّ الَّتِي وَرَسَ فِي مَدْرَسَةٌ لَيْسَتْ فَرَنْسِيّ الَّتِي إِسْتَخْدَمَ نَاسٌ فِي شَوَارِغٌ أَدَانَ .. كُلِّ مَقَامٌ مَقَالٌ

مل احتاج إلى ترجمة كيي تفهم خطاب ملك ؟ لغة " كلاسيكيي " (فصحي) موجود في كل لغة كذلك لغة " دارج ".. فرنسيي التيي درس فيي مدرسة ليست فرنسيي التيي استخدم ناس فيي شوارع باريس .. ملك بريطانيا لا خطب بلغة شوارع أدان .. كل مقام مقال

مل احتاج اليي ترجمه كيي تفهم خطاب ملك ؟ لغه " كلاسيكي "(فصحيي) موجود في كل لغه كذلك لغه " دارج "٠٠ فرنسي التيي درس في مدرسه ليست فرنسي التي استخدم ناس فيي شوارع باريس ٠٠ ملك بريطانيا لا خطب بلغه شوارع ادان ٠٠ كل مقام مقال

هل', 'احتاج', 'إلى', 'ترجمة', 'كين', 'تف', 'خطاب', 'ملك', '؟', 'لغة', '"', 'كلاسيكين', '"(', 'فصحن', ')', 'موجود', 'فين', 'كذلك', 'لغة', '"لنة', 'كذلك', 'الغة', '"ليست', 'فين', 'مدرس', 'فين', 'لغة', 'التين', 'درس', 'فين', 'مدرس', 'لغة', 'التين', 'التين', 'الغة', '"ل

(', 'stopword'), ('"', 'stopword'), ('"', 'stopword'), ('ترجمة', 'noun'), ('"', 'pounct'), ('"', 'pounct'), ('"', 'pounct'), ('"', 'pounct'), ('"', 'pounct'), ('"', 'noun'), ('"', 'noun'), ('"', 'pounct'), ('"', 'pounct'), ('"', 'noun'), ('"', 'stopword'), ('"', 'stopword'), ('"', 'noun'), ('"', 'pounct'), ('"

Stored in directory: c:\users\fady\appdata\local\pip\cache\wheels\ad\e2\59\dbb52181bcb1551e623e55c2a715add67a2fdbca7d62632da8

 $Stored in directory: c: \users fady \appdata \ocal \pip \cache \wheels \10 \46 \58 \dbd 81132 \end{2} e05 \end{2} e2236 \end{2} e2a7 \ff 9d8 \affine \end{2} e769 \end{2} e16998 \end{2} e16999 \end{2} e16999 \end{2} e16999 \end{2} e16999 \end{2} e169999 \end{2} e169999 \end{2} e169999 \end{2} e169999 \end{2} e1699999 \end{2} e169999 \end{2} e169999 \end{2} e169999 \end{2} e1699999 \end{2} e169999 \end{2} e1699999 \end{2} e169999 \end{2} e1699999 \end{2} e169999 \end{2} e169999 \end{2} e169999 \end{2} e169999 \end{2} e169999 \end{2} e16$

['خدم', 'ناس', 'فين', 'شوارع', 'باريس', '..', 'ملك', 'بريطانين', 'لا', 'خطب', 'بلغة', 'شوارع', 'أدان', '..', 'كل', 'مقام', 'مقال

[''اِسْتَخْدَمَ', 'نَاسٌ', 'فِيّ, 'شَوَارِغْ', 'باريس', '..', امَلَكُّ', 'برِيطانِيا', 'لَا', 'خَطَبَ', 'بَلَغَةٌ', 'شَوَارِغْ', 'أَدَانَ', '..', 'كُلِّ', 'مَقَامٌٰ', 'مَقَالٌ

filtered_sentence = []

for word in text.split(" "):

#print(filtered_sentence)

return ps_text , lemmatizer_text

• Hint Word tokenizer: is like split but split depend on space or symbol to seprate words

and Python is awesome. The sky is pinkish-blue. You shouldn't eat cardboard.

In [70]: text = "my name's hossam and i work as a teacher assistant at bfcai and i have 600\$"

s = 'Good muffins cost \$3.88\nin New York. Please buy me two of them.\n\nThanks.'

['Good muffins cost \$3.88\nin New York.', 'Please buy me two of them.', 'Thanks.']

['my', "name's", 'hossam', 'and', 'i', 'work', 'as', 'a', 'teacher', 'assistant', 'at', 'bfcai', 'and', 'i', 'have', '600\$']

['my', 'name', "'s", 'hossam', 'and', 'i', 'work', 'as', 'a', 'teacher', 'assistant', 'at', 'bfcai', 'and', 'i', 'have', '600', '\$']

ا مل تعلم؟ #نيوكاسل يتفوق بالمواجهات المباشرة على #ارسنال في تاريخ الدوري الممتاز الانجليزي؟ بعد قليل أرقام واحمائيات للفريقين'=EXAMPLE_TEXT

['هل', 'تعلم؟', '#', 'نيوكاسل', 'يتفوق', 'بالمواجهات', 'المباشرة', 'على', '#', 'ارسنال', 'في', 'تاريخ', 'الدوري', 'الممتاز', 'الانجليزي؟', 'بعد', 'قليل', 'أرقام', 'واحصائيات', 'للفريقين']

Answer: ["This", "product", "is", "amazing", ",", "but", "the", "delivery", "was", "late", "."]

Tokenization: split paragraph into sentence and called sentence tokenizer

• It contains text processing libraries for tokenization, parsing, classification, stemming, tagging and semantic reasoning

2.2. Punctuation

program : program programs : program programmer : programm programming : program programmers : programm In [22]: 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 HINT: snowballstemmer is somewhat faster and more logical than the original Porter Stemmer. In [23]: **from** nltk.stem **import** SnowballStemmer snowball = SnowballStemmer(language='english') words = ['generous', 'generate', 'generously', 'generation'] for word in words: print(word, "--->", snowball.stem(word)) generous ---> generous generate ---> generat generously ---> generous generation ---> generat Snowball stemmer support different languages In [40]: **from** nltk.stem.snowball **import** SnowballStemmer print(", ".join(SnowballStemmer.languages)) arabic, danish, dutch, english, finnish, french, german, hungarian, italian, norwegian, porter, portuguese, romanian, russian, spanish, swedish stemming in AR In [24]: #stemming is very weak in arabic words = ['ایجری', 'یجرون', 'یجرون', 'عجری' الحریی', 'تجری', 'یجری', الحری for word in words: print (word+' --> '+ ps.stem(word)) الجري --> الجري تجری --> تجری يجرون --> يجرون جری --> جری یجری --> یجری In [25]: **from** nltk.stem.snowball **import** SnowballStemmer s_stemmer = SnowballStemmer(language='arabic') words = ['الجري', 'يجري', 'يجرون', 'جري', 'يجري'] for word in words: print (word+' --> '+s_stemmer.stem(word)) الجري --> الجر تجری --> تجر يجرون --> يجرو *جر*ی --> جر یجری --> یجر In [26]: **from** nltk.stem.isri **import** ISRIStemmer st = ISRIStemmer() words = ['ایجری','یجرون','بیجرون','تجریی', تجریی', for word in words : print(st.stem(word)) print (st.stem(''اعلامیون)) يبجر ملح Lemmatization using nltk Examples of lemmatization: rocks : rock corpora : corpus better : good In [92]: from nltk.stem import WordNetLemmatizer lemmatizer = WordNetLemmatizer() In [94]: print("rocks :", lemmatizer.lemmatize("rocks")) print("corpora :", lemmatizer.lemmatize("corpora")) print("better :", lemmatizer.lemmatize("better", pos ="a")) rocks : rock corpora : corpus better : good In [96]: from nltk.stem import WordNetLemmatizer # Without POS (default is noun) print(lemmatizer.lemmatize("better")) # With POS as adjective ('a') print(lemmatizer.lemmatize("better", pos="a")) # With POS as verb ('v') print(lemmatizer.lemmatize("running", pos="v")) better good run In [98]: # single word lemmatization examples list1 = ['kites', 'babies', 'dogs', 'flying', 'smiling', 'driving', 'died', 'tried', 'feet'] for words in list1: print(words + " ---> " + lemmatizer.lemmatize(words)) kites ---> kite babies ---> baby dogs ---> dog flying ---> flying smiling ---> smiling driving ---> driving died ---> died tried ---> tried feet ---> foot In [30]: words = ["is", "was", "be", "been", "are", "were"] for word in words : print(lemmatizer.lemmatize(word)) is wa be been are were In [31]: #what if we put them as verb words = ["is", "was", "be", "been", "are", "were"] for word in words :