

CSE6060

Statistical Natural Language Processing

NLP - தமிழ்

Name : Kavianand G

Reg. No. : 19MAI0050

Date : 27 – June – 2020

Indic NLP Library

The library provides the following functionalities:

- Text Normalization
- Script Information
- Tokenization
- Word Segmentation
- Script Conversion
- Romanization
- Indicization
- Transliteration
- Translation

In [1]:

```
1 # The path to the local git repo for Indic NLP Library
2 INDIC_NLP_LIB_HOME=r"C:\Users\kavianand\Documents\src\indic_nlp_library"
3
4 # The path to the local git repo for Indic NLP Resources
5 INDIC_NLP_RESOURCES=r"C:\Users\kavianand\Documents\src\indic_nlp_resources"
```

In [2]:

```
1 import sys
2 sys.path.append(r'{}\src'.format(INDIC_NLP_LIB_HOME))
```

In [3]:

```
1 from indicnlp import common
2 common.set_resources_path(INDIC_NLP_RESOURCES)
```

In [4]:

```
1 from indicnlp import loader
2 loader.load()
```

In [5]:

```
1 data="""தமிழ் இலக்கியம் இரண்டாயிரம் ஆண்டுகளுக்கு மேலான தொடர்ச்சி கொண்
2
3 print(data)
```

"தமிழ் இலக்கியம் இரண்டாயிரம் ஆண்டுகளுக்கு மேலான தொடர்ச்சி கொண்ட உலகின் சிறந்த இலக்கியங்களில் ஒன்று. வாழ்வின் பல்வேறு கூறுகளை தமிழ் இலக்கியங்கள் இயம்புகின்றன. தமிழ் இலக்கியத்தில் வெண்பா, குறள், புதுக்கவிதை, கட்டுரை, பழமொழி, தொண்ணூற்றாறு வகை சிற்றிலக்கியங்கள் என பல வடிவங்கள் உள்ளன. தமிழில் வாய்மொழி இலக்கியங்களும் முக்கிய இடம் வகிக்கின்றன.

Text Normalization

In [6]:

```
1 from indicnlp.normalize.indic_normalize import IndicNormalizerFactory
2
3 input_text= data
4
5 factory=IndicNormalizerFactory()
6 normalizer=factory.get_normalizer("ta")
7 output_text=normalizer.normalize(input_text)
8
9 print(input_text)
10 print()
11
12 print('Before normalization')
13 print(' '.join([ c for c in input_text ] ))
14 print('Length: {}'.format(len(input_text)))
15 print()
16 print('After normalization')
17 print(' '.join([ c for c in output_text ] ))
18 print('Length: {}'.format(len(output_text)))
19
```

"தமிழ் இலக்கியம் இரண்டாயிரம் ஆண்டுகளுக்கு மேலான தொடர்ச்சி கொண்ட உலகின் சிறந்த இலக்கியங்களில் ஒன்று. வாழ்வின் பல்வேறு கூறுகளை தமிழ் இலக்கியங்கள் இயம்புகின்றன. தமிழ் இலக்கியத்தில் வெண்பா, குறள், புதுக்கவிதை, கட்டுரை, பழமொழி, தொண்ணூற்றாறு வகை சிற்றிலக்கியங்கள் என பல வடிவங்கள் உள்ளன. தமிழில் வாய்மொழி இலக்கியங்களும் முக்கிய இடம் வகிக்கின்றன.

Before normalization

" த ம ி ழ ி இ ல க ி க ி ய ம ி இ ர ண ஂ ட ா ய ி ர ம ி ஆ ண ஂ ட ஁ க ள ஁ க ி க ி ம ற ி ல ா ன த ி ட ர ி ச ி க ி ண ஂ ட உ ல க ி ன ி ச ி ற ன ி த இ ல க ி க ி ய ன ி க ள ி ல ி ஒ ன ி ற ஁ . வ ா ழ ி வ ி ன ி ப ல ி வ ற ஁ க ி ற ஁ க ள ன ி த ம ி ழ ி இ ல க ி க ி ய ன ி க ள ி இ ய ம ி ப ஁ க ி ன ி ற ன . த ம ி ழ ி இ ல க ி க ி ய த ி த ி ல ி வ ி ண ி ப ா , க ி ற ள ி , ப ஁ த ஁ க ி க வ ி த ன , க ி ட ி ட ஁ ர ன , ப ழ ம ி ழ ி , த ி ண ி ண ி ற ி ற ா ற ஁ வ க ன ி ச ி ற ி ற ி ல க ி க ி ய ன ி க ள ி என ப ல வ ி ல வ ன ி க ள ி உ ள ன ன . த ம ி ழ ி ல ி வ ா ய ி ம ி ழ ி இ ல க ி க ி ய ன ி க ள ி ம ி ம ஁ க ி க ி ய இ ட ம ி வ க ி க ி க ி ன ி ற ன .

Length: 338

After normalization

" த ம ி ழ ி இ ல க ி க ி ய ம ி இ ர ண ஂ ட ா ய ி ர ம ி ஆ ண ஂ ட ஁ க ள ஁ க ி க ி ம ற ி ல ா ன த ி ட ர ி ச ி க ி ண ஂ ட உ ல க ி ன ி ச ி ற ன ி த இ ல க ி க ி ய ன ி க ள ி ல ி ஒ ன ி ற ஁ . வ ா ழ ி வ ி ன ி ப ல ி வ ற ஁ க ி ற ஁ க ள ன ி த ம ி ழ ி இ ல க ி க ி ய ன ி க ள ி இ ய ம ி ப ஁ க ி ன ி ற ன . த ம ி ழ ி இ ல க ி க ி ய த ி த ி ல ி வ ி ண ி ப ா , க ி ற ள ி , ப ஁ த ஁ க ி க வ ி த ன , க ி ட ி ட ஁ ர ன , ப ழ ம ி ழ ி , த ி ண ி ண ி ற ி ற ா ற ஁ வ க ன ி ச ி ற ி ற ி ல க ி க ி ய ன ி க ள ி என ப ல வ ி ல வ ன ி க ள ி உ ள ன ன . த ம ி ழ ி ல ி வ ா ய ி ம ி ழ ி இ ல க ி க ி ய ன ி க ள ி ம ி ம ஁ க ி க ி ய இ ட ம ி வ க ி க ி க ி ன ி ற ன .

Length: 338

Sentence Tokenization / Sentence Splitter

In [7]:

```
1 from indicnlp.tokenize import sentence_tokenize
2
3 string = data
4 sentences=sentence_tokenize.sentence_split(string, lang='ta')
5 for t in sentences:
6     print("\n",t)
```

"தமிழ் இலக்கியம் இரண்டாயிரம் ஆண்டுகளுக்கு மேலான தொடர்ச்சி கொண்ட உலகின் சிறந்த இலக்கியங்களில் ஒன்று.

வாழ்வின் பல்வேறு கூறுகளை தமிழ் இலக்கியங்கள் இயம்புகின்றன.

தமிழ் இலக்கியத்தில் வெண்பா, குறள், புதுக்கவிதை, கட்டுரை, பழமொழி, தொண்ணூற்றாறு வகை சிற்றிலக்கியங்கள் என பல வடிவங்கள் உள்ளன.

தமிழில் வாய்மொழி இலக்கியங்களும் முக்கிய இடம் வகிக்கின்றன.

Word Tokenization

In [8]:

```
1 from indicnlp.tokenize import indic_tokenize
2 from indicnlp.tokenize import indic_detokenize
3
4 string=data
5
6 print('\nInput String: {}'.format(string))
7 print('\nTokens: ')
8 for t in indic_tokenize.trivial_tokenize(string):
9     print(t)
```

Input String: "தமிழ் இலக்கியம் இரண்டாயிரம் ஆண்டுகளுக்கு மேலான தொடர்ச்சி கொண்ட உலகின் சிறந்த இலக்கியங்களில் ஒன்று. வாழ்வின் பல்வேறு கூறுகளை தமிழ் இலக்கியங்கள் இயம்புகின்றன. தமிழ் இலக்கியத்தில் வெண்பா, குறள், புதுக்கவிதை, கட்டுரை, பழமொழி, தொண்ணூற்றாறு வகை சிற்றிலக்கியங்கள் என பல வடிவங்கள் உள்ளன. தமிழில் வாய்மொழி இலக்கியங்களும் முக்கிய இடம் வகிக்கின்றன."

Tokens:

```
"
தமிழ்
இலக்கியம்
இரண்டாயிரம்
ஆண்டுகளுக்கு
மேலான
தொடர்ச்சி
கொண்ட
உலகின்
சிறந்த
இலக்கியங்களில்
ஒன்று
.
வாழ்வின்
பல்வேறு
கூறுகளை
தமிழ்
இலக்கியங்கள்
இயம்புகின்றன
.
தமிழ்
இலக்கியத்தில்
வெண்பா
,
குறள்
,
புதுக்கவிதை
,
கட்டுரை
,
பழமொழி
,
தொண்ணூற்றாறு
வகை
சிற்றிலக்கியங்கள்
என
பல
வடிவங்கள்
உள்ளன
.
```

தமிழில்
வாய்மொழி
இலக்கியங்களும்
முக்கிய
இடம்
வகிக்கின்றன
.

Syllabification

In [9]:

```
1 from indicnlp.syllable import syllabifier
2
3 text = 'சிட்டுக்குருவி மனித குடியிருப்புடன் வலுவாகத் தொடர்புடையது ஆகுமீ'
4 lang='ta'
5
6 syl_word = (' '.join(syllabifier.orthographic_syllabify(text,lang)))
7 print("Syllabification (Approximation)\n")
8 print(syl_word)## Transliteration
```

Syllabification (Approximation)

சி ட்டு க் கு ரு வி ம னி த கு டி யி ரு ப் பு ட ன் வ லு வாக த் தொ
ட ர்பு டை ய து ஆ கு மீ

Morphological Analyser

In [10]:

```
1 from indicnlp.morph import unsupervised_morph
2 from indicnlp import common
3
4 analyzer=unsupervised_morph.UnsupervisedMorphAnalyzer('ta')## Syllabification
```

In [11]:

```
1 string = 'புத்துணர்ச்சியான'
2
3 analyzes_tokens = analyzer.morph_analyze_document(string.split(' '))
4 print(analyzes_tokens)
```

['புத்துணர்ச்சி', 'யான']

In [12]:

```
1 text = 'சிட்டுக்குருவி மனித குடியிருப்புடன் வலுவாகத் தொடர்புடையது ஆகும்'
2
3 analyzes_tokens=analyzer.morph_analyze_document(text.split(' '))
4
5 for w in analyzes_tokens:
6     print(w)
```

சிட்டுக்குருவி
மனித
குடியிருப்பு
டன்
வலுவாக
த்
தொடர்புடைய
து
ஆகும்
ி

POS Tagging - தமிழ்

In [13]:

```
1 from rippletagger.tagger import Tagger
2
3 def pos_tag_tamil(query):
4     tagger = Tagger(language="tam")
5     posTagger = tagger.tag(query)
6     return(posTagger)
```

In [14]:

```
1 query = 'சிட்டுக்குருவி மனித குடியிருப்புடன் வலுவாகத் தொடர்புடையது ஆகும்'
2 tamil_tag = pos_tag_tamil(query)
3 print(tamil_tag)
```

[('சிட்டுக்குருவி', 'NOUN'), ('மனித', 'PROPN'), ('குடியிருப்புடன்', 'NOUN'), ('வலுவாகத்', 'PROPN'), ('தொடர்புடையது', 'VERB'), ('ஆகும்', 'VERB')]

In [15]:

```
1 query = "தமிழ் இலக்கியம் இரண்டாயிரம் ஆண்டுகளுக்கு மேலான தொடர்ச்சி கொண்  
2 tamil_tag = pos_tag_tamil(query)  
3 print(tamil_tag)
```

```
[('தமிழ்', 'PROPN'), ('இலக்கியம்', 'NOUN'), ('இரண்டாயிரம்', 'NOUN'),  
( 'ஆண்டுகளுக்கு', 'NOUN'), ('மேலான', 'ADJ'), ('தொடர்ச்சி', 'NOUN'),  
( 'கொண்ட', 'ADP'), ('உலகின்', 'NOUN'), ('சிறந்த', 'ADJ'), ('இலக்கிய  
ங்களில்', 'NOUN'), ('ஒன்று.', 'NOUN'), ('வாழ்வின்', 'NOUN'), ('பல்வேறு',  
'ADJ'), ('கூறுகளை', 'NOUN'), ('தமிழ்', 'PROPN'), ('இலக்கியங்கள்', 'NOUN'),  
( 'இயம்புகின்றன.', 'NOUN'), ('தமிழ்', 'PROPN'), ('இலக்கியத்தில்',  
'NOUN'), ('வெண்பா,', 'NOUN'), ('குறள்,', 'NOUN'), ('புதுக்கவிதை,', 'NOUN'),  
( 'கட்டுரை,', 'NOUN'), ('பழமொழி,', 'NOUN'), ('தொண்ணூற்றாறு',  
'NOUN'), ('வகை', 'NOUN'), ('சிறிறிலக்கியங்கள்', 'NOUN'), ('என', 'PART'),  
( 'பல', 'ADJ'), ('வடிவங்கள்', 'NOUN'), ('உள்ளன.', 'NOUN'), ('த  
மிழில்', 'NOUN'), ('வாய்மொழி', 'NOUN'), ('இலக்கியங்களும்', 'ADJ'),  
( 'முக்கிய', 'ADJ'), ('இடம்', 'ADP'), ('வகிக்கின்றன.', 'VERB')]
```

StopWords - தமிழ்

In [16]:

```
1 with open('TamilStopWords.txt', encoding="utf8") as file:  
2     stop_words_tamil = file.read()  
3  
4 print(stop_words_tamil)
```

ஒரு
என்று
மற்றும்
இந்த
இது
என்ற
கொண்டு
என்பது
பல
ஆகும்
அல்லது
அவர்
நான்
உள்ள
அந்த
இவர்
என
முதல்
என்ன
...-...-...

In [17]:

```
1 data=""  
2
```


In [18]:

```
1 token = []
2 for t in indic_tokenize.trivial_tokenize(data):
3     token.append(t)
```

In [19]:

```
1 print("Tokenized Words :\n ")
2 print(token)
```

Tokenized Words :

['', 'தமிழ்', 'இலக்கியம்', 'இரண்டாயிரம்', 'ஆண்டுகளுக்கு', 'மேலான', 'தொடர்ச்சி', 'கொண்ட', 'உலகின்', 'சிறந்த', 'இலக்கியங்களில்', 'ஒன்று', '.', 'வாழ்வின்', 'பல்வேறு', 'கூறுகளை', 'தமிழ்', 'இலக்கியங்கள்', 'இயம்புகின்றன', '.', 'தமிழ்', 'இலக்கியத்தில்', 'வெண்பா', ',', 'குறள்', ',', 'புதுக்கவிதை', ',', 'கட்டுரை', ',', 'பழமொழி', ',', 'தொண்ணூற்றாறு', 'வகை', 'சிறுநிலக்கியங்கள்', 'என', 'பல', 'வடிவங்கள்', 'உள்ளன', '.', 'தமிழில்', 'வாய்மொழி', 'இலக்கியங்களும்', 'முக்கிய', 'இடம்', 'வகிக்கின்றன', '.']

In [20]:

```
1 fil_sen = []
2 fil_sw = []
3
4 for t in indic_tokenize.trivial_tokenize(data):
5     if t not in stop_words_tamil:
6         fil_sen.append(t)
7     else:
8         fil_sw.append(t)
```

In [21]:

```
1 print(fil_sen)
```

['', 'தமிழ்', 'இலக்கியம்', 'இரண்டாயிரம்', 'ஆண்டுகளுக்கு', 'மேலான', 'தொடர்ச்சி', 'உலகின்', 'சிறந்த', 'இலக்கியங்களில்', 'ஒன்று', '.', 'வாழ்வின்', 'கூறுகளை', 'தமிழ்', 'இலக்கியங்கள்', 'இயம்புகின்றன', '.', 'தமிழ்', 'இலக்கியத்தில்', 'வெண்பா', ',', 'குறள்', ',', 'புதுக்கவிதை', ',', 'கட்டுரை', ',', 'பழமொழி', ',', 'தொண்ணூற்றாறு', 'வகை', 'சிறுநிலக்கியங்கள்', 'வடிவங்கள்', '.', 'தமிழில்', 'வாய்மொழி', 'இலக்கியங்களும்', 'முக்கிய', 'வகிக்கின்றன', '.']

In [22]:

```
1 print("Removed Stop Words :\n ")
2 print(fil_sw)
```

Removed Stop Words :

['கொண்ட', 'பல்வேறு', 'என', 'பல', 'உள்ளன', 'இடம்']

In [23]:

```
1 print("Number of Words including StopWords : " ,len(token))
2 print("Number of Words excluding StopWords : " ,len(fil_sen))
3 print("Total number of Stop Words Excluded : " , len(fil_sw))
```

Number of Words including StopWords : 47

Number of Words excluding StopWords : 41

Total number of Stop Words Excluded : 6

Transliteration

In [24]:

```
1 from indicnlp.transliterate.unicode_transliterate import UnicodeIndicTransliterator
2
3 input_text = 'சிட்டுக்குருவி மனித குடியிருப்புடன் வலுவாகத் தொடர்புடையது ஆகும்'
4 ta_hi_text = (UnicodeIndicTransliterator.transliterate(input_text,"ta","hi"))
5
6 print("Transliterate From Tamil to Hindi\n")
7 print("Input Text : " , input_text)
8 print("output Text : ", ta_hi_text)
```

Transliterate From Tamil to Hindi

Input Text : சிட்டுக்குருவி மனித குடியிருப்புடன் வலுவாகத் தொடர்புடையது ஆகும்

output Text : चिट्ठुक्कुरुवि मनिथ कुटियिरुप्पुटन् वलुवाकत् तोटपुटैयतु आकुम्

In [25]:

```
1 from indicnlp.transliterate.unicode_transliterate import UnicodeIndicTransliterator
2
3 input_text=ta_hi_text
4 hi_ta_text = (UnicodeIndicTransliterator.transliterate(input_text,"hi","ta"))
5
6 print("Transliterate From Hindi to Tamil\n")
7 print("Input Text : " , input_text)
8 print("output Text : ", hi_ta_text)
```

Transliterate From Hindi to Tamil

Input Text : चिट्ठुक्कुरुवि मनिथ कुटियिरुप्पुटन् वलुवाकत् तोटपुटैयतु आकुम्

output Text : சிட்டுக்குருவி மனித குடியிருப்புடன் வலுவாகத் தொடர்புடையது ஆகும்

In [26]:

```
1 from indicnlp.transliterate.unicode_transliterate import ItransTransliterator
2
3 input_text = 'சிட்டுக்குருவி மனித குடியிருப்புடன் வலுவாகத் தொடர்புடையது ஆகும்'
4 lang='ta'
5
6 print(ItransTransliterator.to_itrans(input_text,lang))
7
8 print("Transliterate From Hindi to Tamil\n")
9 print("Input Text : " , input_text)
10 print("output Text : ", hi_ta_text)
```

chiTTukkuruvi ma*nita kuTiyiruppuTa*n valuvaakat t.oTarpuTaiyatu aakum
Transliterate From Hindi to Tamil

Input Text : சிட்டுக்குருவி மனித குடியிருப்புடன் வலுவாகத் தொடர்புடையது ஆகும்
output Text : சிட்டுக்குருவி மனித குடியிருப்புடன் வலுவாகத் தொடர்புடையது ஆகும்

Translation Options

In [27]:

```
1 text ='சிட்டுக்குருவி மனித குடியிருப்புடன் வலுவாகத் தொடர்புடையது ஆகும்'## St
```

In [28]:

```
1 from textblob import TextBlob
2 txt = TextBlob(text)
3 # Can detect language
4 txt.detect_language()
```

Out[28]:

'ta'

In [29]:

```
1 # Language Translation also available in TextBlob
2 ta_eng_text = txt.translate(from_lang='ta', to ='en')
3 print(ta_eng_text)
```

The sparrow is strongly associated with human habitation

In [30]:

```
1 txt = ta_eng_text
2 # Language Translation also available in TextBlob
3 print(txt.translate(from_lang='en', to ='ta'))
```

குருவி மனித வாழ்விடத்துடன் வலுவாக தொடர்புடையது

