

# NLTK 자연어 처리 패키지

1692047 원진

# 설치

## Mac / Unix

1. install NLTK
2. install Numpy (optional)
3. test installation

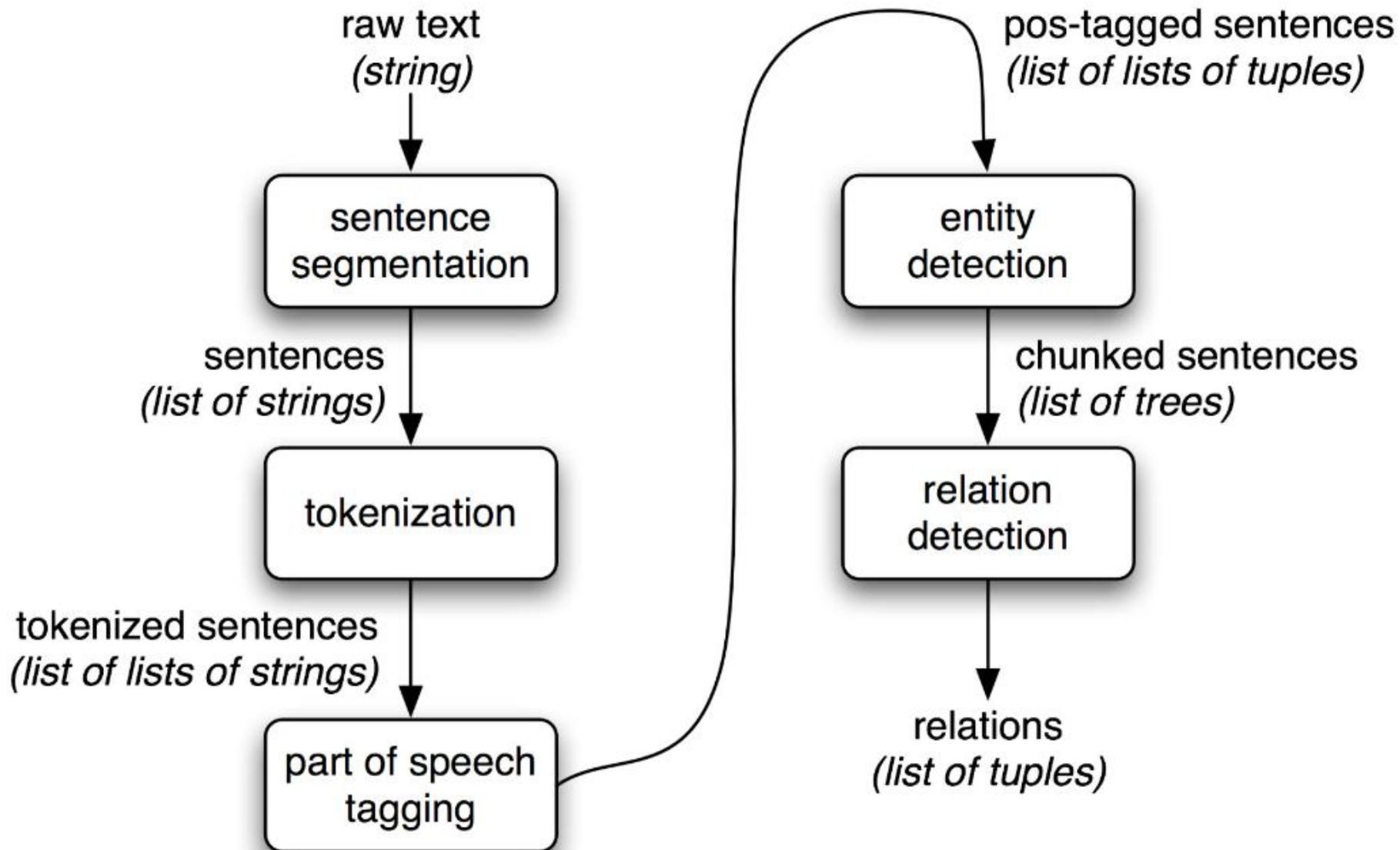
```
run sudo pip install -U nltk  
run sudo pip install -U numpy  
run python then type import nltk
```

## Windows

1. install Python 3.7
2. install Numpy (optional)
3. install NLTK
4. test installation

```
http://www.python.org/downloads/  
(avoid the 64-bit versions)  
https://www.scipy.org/scipylib/download.html  
http://pypi.python.org/pypi/nltk  
start->Python37, then type import nltk
```

# Information Extraction Architecture



# 토큰 생성

**token** 긴 문자열을 분석을 위한 작은 단위로 나눈 것.

**tokenizer** 문자열을 토큰으로 분리하는 함수, 토큰 문자열의 리스트를 반환한다.

- `sent_tokenize(text, language = 'english')`

```
1 import nltk
2
3 corpus = 'hello, my name is jin! nice to meet you.'
4 print(nltk.sent_tokenize(corpus))
```

```
['hello, my name is wonjin!', 'nice to meet you.']
```

- `word_tokenize(text, language='english', preserve_line=False)`

```
1 import nltk
2
3 corpus = 'hello, my name is jin! nice to meet you.'
4 print(nltk.word_tokenize(corpus))
```

```
['hello', ',', 'my', 'name', 'is', 'jin', '!', 'nice', 'to', 'meet', 'you', '.']
```

# 토큰 생성

- 영문의 경우 정규 표현식을 쓸 수 있다.

```
1 import nltk
2
3 corpus = 'hello, my name is jin! nice to meet you.'
4
5 retokenize = nltk.tokenize.RegexpTokenizer("[\w]+")
6 retokenize.tokenize(corpus)
7
8 print(retokenize.tokenize(corpus))
```

```
['hello', 'my', 'name', 'is', 'jin', 'nice', 'to', 'meet', 'you']
```

# 형태소 분석

## 형태소 (morpheme)

언어학에서 일정한 의미가 있는 가장 작은 말의 단위.  
보통 자연어 처리에서는 토큰으로 형태소를 이용함.

## 형태소 분석 (morphological analysis)

- 어간 추출 (stemming)
- 원형 복원 (lemmatizing)
- 품사 부착 (Part-Of-Speech tagging)

단어로부터 어근, 접두사, 접미사, 품사 등  
다양한 언어적 속성을 파악하고,  
이를 이용하여  
형태소를 찾아내거나 처리하는 작업.

# 형태소 분석 - 어간 추출 (stemming)

변화된 단어의 접미사나 어미를 제거하여 같은 의미를 가지는 형태소의 기본형을 찾는 방법.  
원형을 정확히 찾아주지는 않는다.

- PorterStemmer *A word stemmer based on the Porter stemming algorithm.*

```
1 import nltk
2
3 words = ['lives', 'crying', 'flies', 'dying']
4
5 st = nltk.stem.PorterStemmer()
6 print([st.stem(w) for w in words])
```

```
['live', 'cri', 'fli', 'die']
```

- LancasterStemmer *A word stemmer based on the Lancaster (Paice/Husk) stemming algorithm.*

```
1 import nltk
2
3 words = ['lives', 'crying', 'flies', 'dying']
4
5 st = nltk.stem.LancasterStemmer()
6 print([st.stem(w) for w in words])
```

```
['liv', 'cry', 'fli', 'dying']
```

# 형태소 분석 - 원형 복원 (lemmatizing)

같은 의미를 가지는 여러 단어를 사전형으로 통일하는 작업.  
품사(part of speech)를 지정하는 경우 좀 더 정확한 원형을 찾을 수 있다.

wordnet 패키지를 추가로 다운받아야 함.

```
>>> import nltk
>>> nltk.download('wordnet')
```

```
1 import nltk
2
3 words = ['lives', 'crying', 'flies', 'dying']
4
5 lm = nltk.stem.WordNetLemmatizer()
6 print([lm.lemmatize(w) for w in words])
```

```
['life', 'cry', 'fly', 'dying']
```

```
1 import nltk
2
3 words = ['lives', 'crying', 'flies', 'dying']
4
5 lm = nltk.stem.WordNetLemmatizer()
6 print(lm.lemmatize("dying", pos="v"))
```

```
die
```



# 형태소 분석 - 품사 부착 (Part-Of-Speech tagging)

낱말을 문법적인 기능이나 형태, 뜻에 따라 구분한 것.  
NLTK에서는 Penn Treebank Tagset을 이용함.

TAG	DESCRIPTION	EXAMPLE
CC	conjunction, coordinating	<i>and, or, but</i>
CD	cardinal number	<i>five, three, 13%</i>
DT	determiner	<i>the, a, these</i>
EX	existential there	<i><u>there</u> were six boys</i>
FW	foreign word	<i>mais</i>
IN	conjunction, subordinating or preposition	<i>of, on, before, unless</i>
JJ	adjective	<i>nice, easy</i>
JJR	adjective, comparative	<i>nicer, easier</i>
JJS	adjective, superlative	<i>nicest, easiest</i>
LS	list item marker	
MD	verb, modal auxillary	<i>may, should</i>

# 형태소 분석 - 품사 부착 (Part-Of-Speech tagging)

TAG	DESCRIPTION	EXAMPLE
NN	noun, singular or mass	<i>tiger, chair, laughter</i>
NNS	noun, plural	<i>tigers, chairs, insects</i>
NNP	noun, proper singular	<i>Germany, God, Alice</i>
NNPS	noun, proper plural	<i>we met two <u>Christmases</u> ago</i>
PDT	predeterminer	<i><u>both</u> his children</i>
POS	possessive ending	<i>'s</i>
PRP	pronoun, personal	<i>me, you, it</i>
PRP\$	pronoun, possessive	<i>my, your, our</i>
RB	adverb	<i>extremely, loudly, hard</i>
RBR	adverb, comparative	<i>better</i>
RBS	adverb, superlative	<i>best</i>
RP	adverb, particle	<i>about, off, up</i>
SYM	symbol	<i>%</i>

# 형태소 분석 - 품사 부착 (Part-Of-Speech tagging)

TAG	DESCRIPTION	EXAMPLE
TO	infinitival to	<i>what <u>to</u> do?</i>
UH	interjection	<i>oh, oops, gosh</i>
VB	verb, base form	<i>think</i>
VBZ	verb, 3rd person singular present	<i>she <u>thinks</u></i>
VBP	verb, non-3rd person singular present	<i>I <u>think</u></i>
VBD	verb, past tense	<i>they <u>thought</u></i>
VCN	verb, past participle	<i>a <u>sunken</u> ship</i>
VBG	verb, gerund or present participle	<i><u>thinking</u> is fun</i>
WDT	<i>wh</i> -determiner	<i>which, whatever, whichever</i>
WP	<i>wh</i> -pronoun, personal	<i>what, who, whom</i>
WP\$	<i>wh</i> -pronoun, possessive	<i>whose, whosever</i>
WRB	<i>wh</i> -adverb	<i>where, when</i>
.	punctuation mark, sentence closer	<i>. ; ? *</i>

# 형태소 분석 - 품사 부착 (Part-Of-Speech tagging)

TAG	DESCRIPTION	EXAMPLE
,	punctuation mark, comma	,
:	punctuation mark, colon	:
(	contextual separator, left paren	(
)	contextual separator, right paren	)

- `pos_tag` 단어 토큰에 품사를 부착하여 튜플로 출력함.  
averaged\_perceptron\_tagger 패키지를 추가로 다운받아야 함.

```
>>> import nltk
>>> nltk.download('averaged_perceptron_tagger')
```

```
1 import nltk
2
3 sentence = "Emma refused to permit us to obtain the refuse permit"
4 tagged_list = nltk.tag.pos_tag(nltk.word_tokenize(sentence))
5
6 print(tagged_list)
```

```
[('Emma', 'NNP'), ('refused', 'VBD'), ('to', 'TO'), ('permit', 'VB'), ('us', 'PRP'), ('to', 'TO'), ('obtain', 'VB'),
```

# 형태소 분석 - 품사 부착 (Part-Of-Speech tagging)

- 품사 태그 정보를 사용하면 명사인 토큰만 선택할 수 있다.

```
1 import nltk
2
3 sentence = "Emma refused to permit us to obtain the refuse permit"
4 tagged_list = nltk.tag.pos_tag(nltk.word_tokenize(sentence))
5
6 nouns_list = [t[0] for t in tagged_list if t[1] == "NN"]
7 print(nouns_list)
```

```
['refuse', 'permit']
```

# 어휘 분석 (Lexical Analysis)

## Part of speech:

NP NNP RB VBD IN NNP NNP , CC PRP VBZ RB VBG PRP IN PRP .  
Mrs. Clinton previously worked for Mr. Obama, but she is now distancing herself from him .

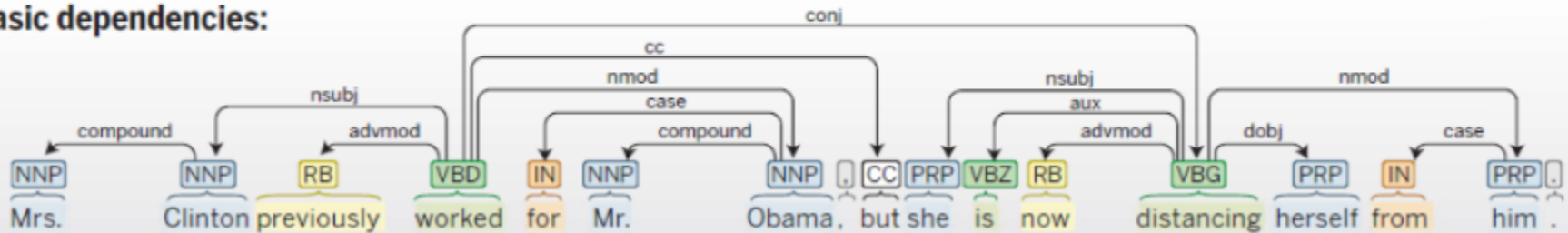
## Named entity recognition:

Person Date Person Date  
Mrs. Clinton previously worked for Mr. Obama, but she is now distancing herself from him.

## Co-reference:

Coref Coref Coref  
Mention Ment M Mention M  
Mrs. Clinton previously worked for Mr. Obama, but she is now distancing herself from him.

## Basic dependencies:



# Named Entity Recognition

## 개체 명 인식 (NER)

인명, 지명 등 고유 명사를 분류하는 방법론.

NE Type	Examples
ORGANIZATION	<i>Georgia-Pacific Corp., WHO</i>
PERSON	<i>Eddy Bonte, President Obama</i>
LOCATION	<i>Murray River, Mount Everest</i>
DATE	<i>June, 2008-06-29</i>
TIME	<i>two fifty a m, 1:30 p.m.</i>
MONEY	<i>175 million Canadian Dollars, GBP 10.40</i>
PERCENT	<i>twenty pct, 18.75 %</i>
FACILITY	<i>Washington Monument, Stonehenge</i>
GPE	<i>South East Asia, Midlothian</i>

NLTK에서 지원하는 기본 모듈로는 정확한 결과를 얻을 수 없다.  
따라서 추가적인 데이터 셋을 확보해야 하며,  
일반적으로 Stanford NER tagger가 사용된다.

# Named Entity Recognition

## Stanford NER Tagger

- Java JRE 8 이상 설치.
- JAVAHOME 환경변수 지정.
- dataset 경로 지정.

```
import nltk
import os
from nltk.tag.stanford import StanfordNERTagger

os.environ['JAVAHOME'] = "C:/Program Files/Java/jre1.8.0_191/bin/java.exe"

jar = './stanford-ner-3.9.2.jar'
model = './english.all.3class.distsim.crf.ser'

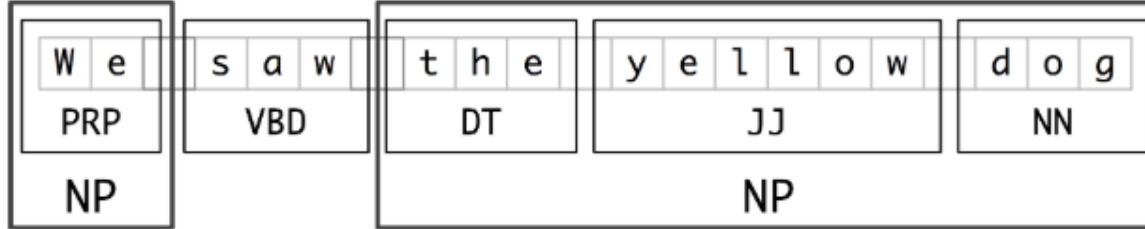
# Prepare NER tagger with english model
ner_tagger = StanfordNERTagger(model, jar, encoding='utf8')
```



# Chunking

## Chunk

구(phrase)로 묶어서 레이블을 지정함.



TAG	DESCRIPTION	WORDS	EXAMPLE
NP	noun phrase	DT+RB+JJ+NN + PR	<i>the strange bird</i>
PP	prepositional phrase	TO+IN	<i>in between</i>
VP	verb phrase	RB+MD+VB	<i>was looking</i>
ADVP	adverb phrase	RB	<i>also</i>
ADJP	adjective phrase	CC+RB+JJ	<i>warm and cosy</i>
SBAR	subordinating conjunction	IN	<i><u>whether</u> or not</i>
PRT	particle	RP	<i><u>up</u> the stairs</i>
INTJ	interjection	UH	<i>hello</i>