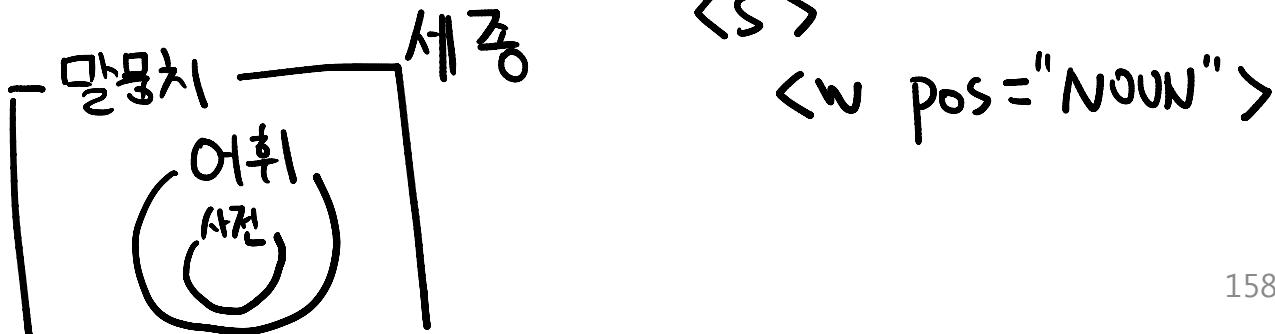
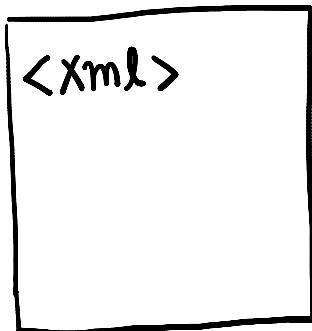
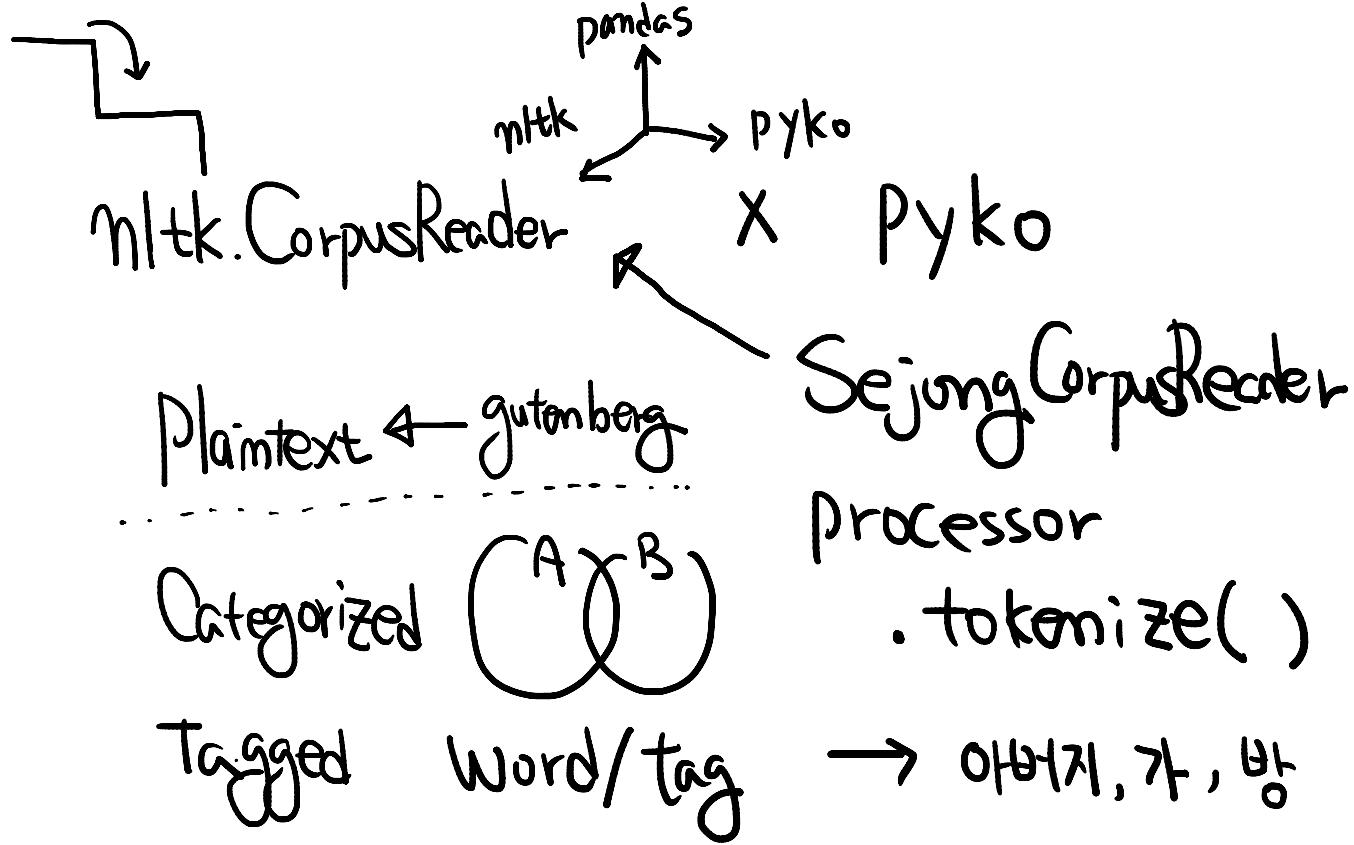
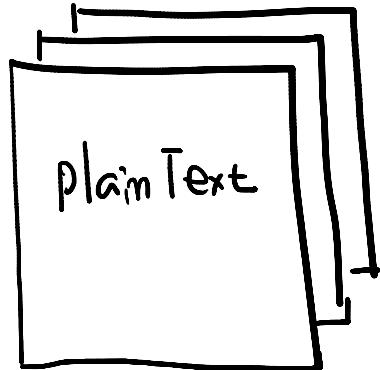


파이썬 자연어 처리

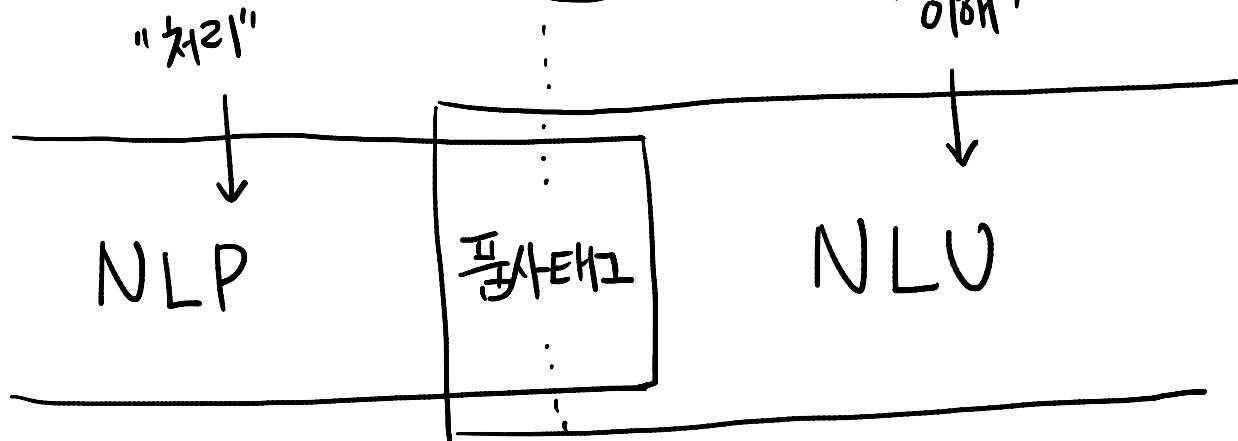
이성주

seongjoo@codebasic.io

Day 3
말뭉치



자연어



말뭉치, 토큰화, ...

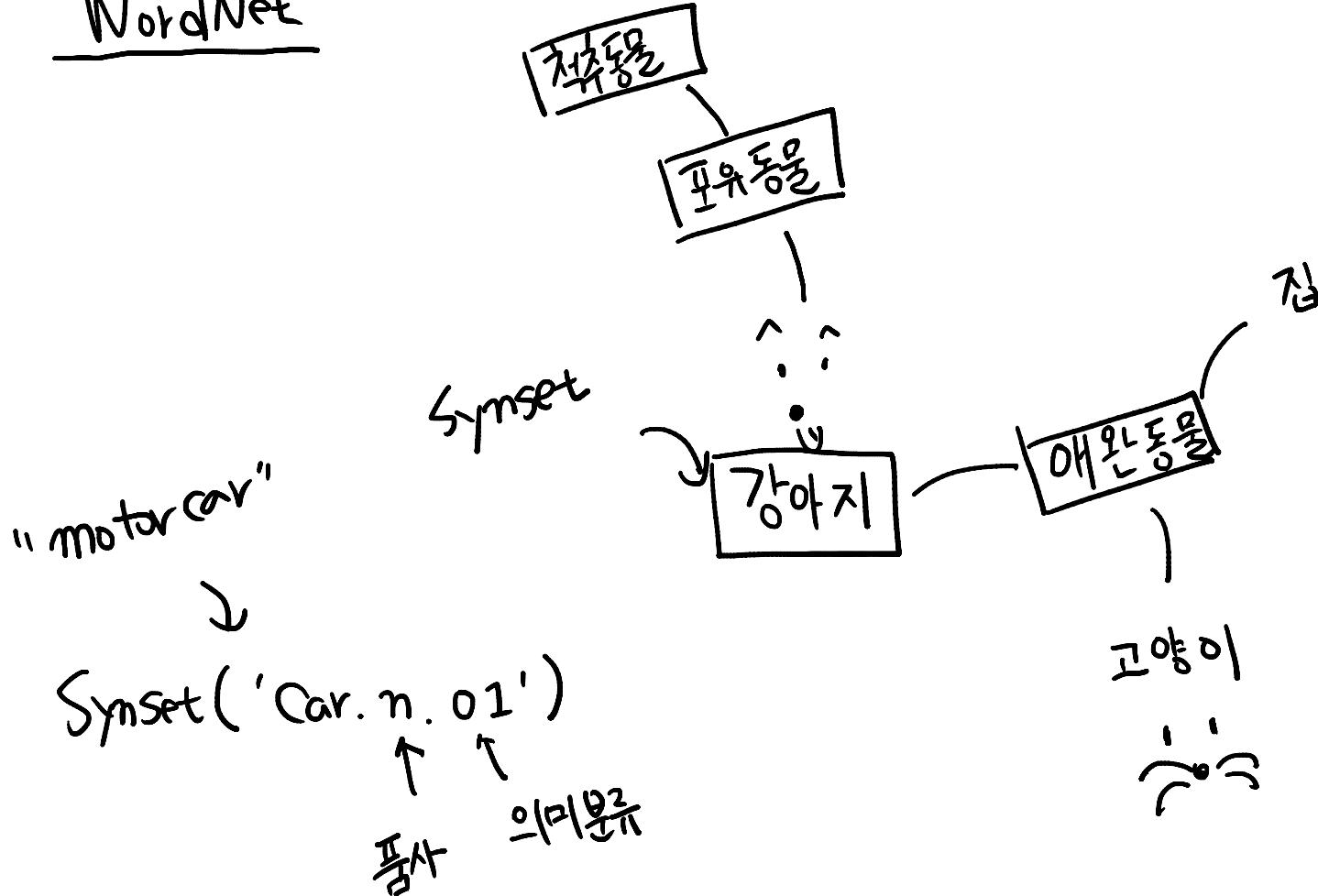
"오늘 날씨가 어때?"

↳ 데이터 선택 ~



품사태그

WordNet



File Edit View Insert Cell Kernel Widgets Help

In [5]: synset = wordnet.synset('car.n.01')

.synsets(정식기配偶)

In [6]: synset.definition()

색인

Out[6]: 'a motor vehicle with four wheels; usually propelled by an internal combustion engine'

In [7]: synset.examples()

상위

Out[7]: ['he needs a car to get to work']

In [8]: synset.lemmas() 동의어

Out[8]: [Lemma('car.n.01.car'), ...
Lemma('car.n.01.auto'),
Lemma('car.n.01.automobile'),
Lemma('car.n.01.machine'),
Lemma('car.n.01.motorcar')]

Synset

Synset

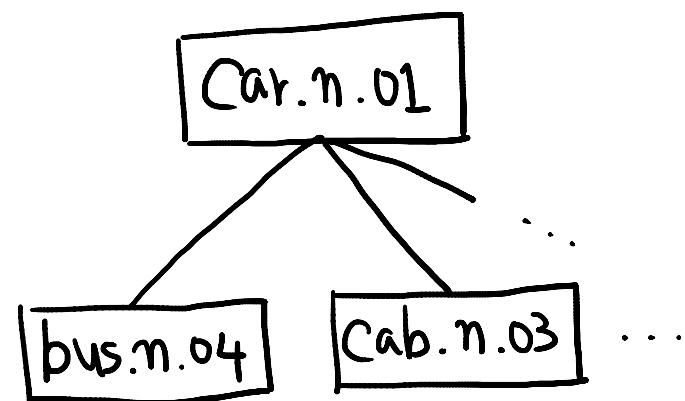
Synset

In []:

```
In [11]: motorcar = wordnet.synset('car.n.01')
```

```
In [12]: motorcar.hyponyms() "하위어"
```

```
Out[12]: [Synset('ambulance.n.01'),  
          Synset('beach_wagon.n.01'),  
          Synset('bus.n.04'),  
          Synset('cab.n.03'),  
          Synset('compact.n.03'),  
          Synset('convertible.n.01'),  
          Synset('coupe.n.01'),  
          Synset('cruiser.n.01'),  
          Synset('electric.n.01'),  
          Synset('gas_guzzler.n.01'),  
          Synset('hardtop.n.01'),  
          Synset('hatchback.n.01'),  
          Synset('horseless_carriage.n.01'),  
          Synset('hot_rod.n.01'),  
          Synset('jeep.n.01'),  
          Synset('limousine.n.01'),  
          Synset('loaner.n.02'),  
          Synset('minicar.n.01'),  
          Synset('minivan.n.01'),  
          Synset('model_t.n.01'),  
          Synset('pace_car.n.01'),
```

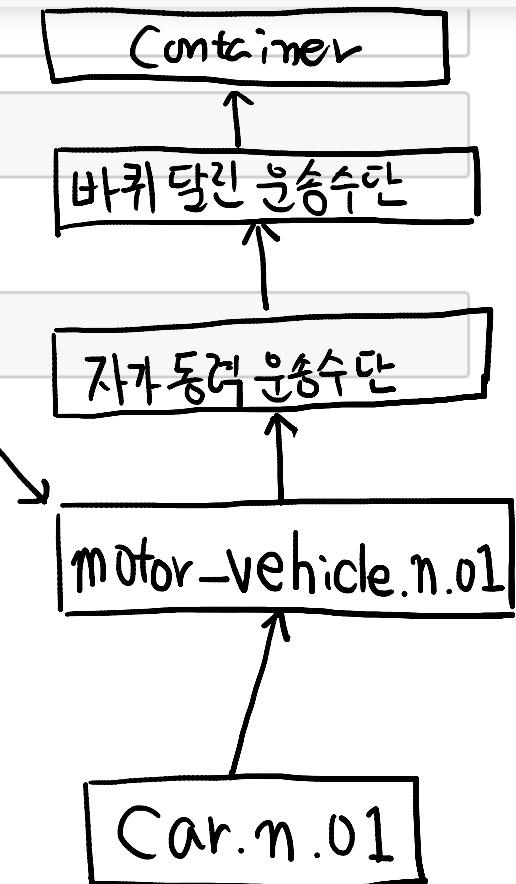


In [13]: motorcar.hypernyms() "상위어"

Out[13]: [Synset('motor_vehicle.n.01')]

In [14]: motorcar.hypernym_paths() "상위어 경로"

Out[14]: ① [Synset('entity.n.01'),
 Synset('physical_entity.n.01'),
 Synset('object.n.01'),
 Synset('whole.n.02'),
 Synset('artifact.n.01'),
 Synset('instrumentality.n.03'),
 Synset('container.n.01'),
 Synset('wheeled_vehicle.n.01'),
 Synset('self-propelled_vehicle.n.01'),
 Synset('motor_vehicle.n.01'),
 Synset('car.n.01')],
 ② [Synset('entity.n.01'),
 Synset('physical_entity.n.01'),
 Synset('object.n.01'),
 Synset('whole.n.02'),
 Synset('artifact.n.01'),
 Synset('instrumentality.n.03'),
 Synset('conveyance.n.03'),
 Synset('vehicle.n.01')]



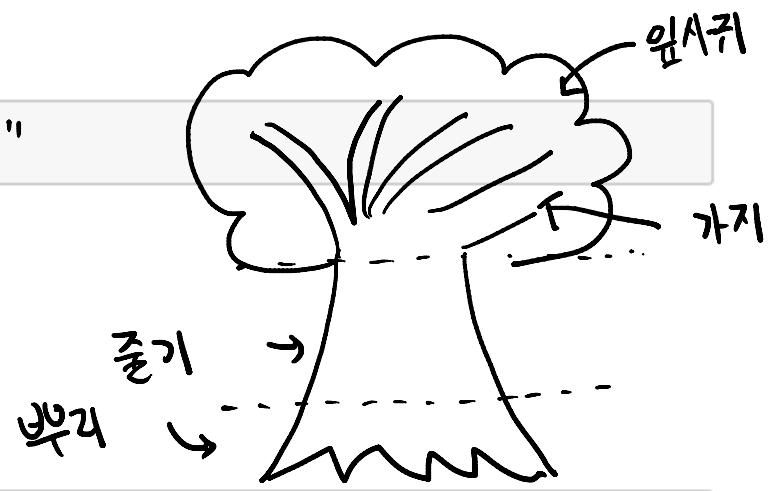
```
In [15]: tree = wordnet.synset('tree.n.01')
```

```
In [16]: tree.definition()
```

Out[16]: 'a tall perennial woody plant having a main trunk and branches forming a distinct elevated crown; includes both gymnosperms and angiosperms'

```
In [17]: tree.part_meronyms() "부분어"
```

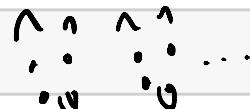
Out[17]: [Synset('burl.n.02'),
Synset('crown.n.07'),
Synset('limb.n.02'),
Synset('stump.n.01'),
Synset('trunk.n.01')]



```
In [18]: tree.member_holonyms() "전체어"
```

Out[18]: [Synset('forest.n.01')]

```
In [ ]:
```



File Edit View Insert Cell Kernel Widgets Help

In [23]: 강아지 = wordnet.synset('dog.n.01')
강아지.definition()

Out[23]: 'a member of the genus Canis (probably descended from the common wolf) that has been domesticated by man since prehistoric times; occurs in many breeds'

In [26]: 고양이 = wordnet.synset('cat.n.01')

동물

In [35]: 사자 = wordnet.synset('lion.n.01')

In [27]: wordnet.synsets('human')

Out[27]: [Synset('homo.n.02'),
Synset('human.a.01'),
Synset('human.a.02'),
Synset('human.a.03')]



In [32]: 인간 = wordnet.synset('homo.n.02')

In []:

File Edit View Insert Cell Kernel Widgets Help

```
Out[37]: [Synset('homo.n.02'),  
          Synset('human.a.01'),  
          Synset('human.a.02'),  
          Synset('human.a.03')]
```

```
In [32]: 인간 = wordnet.synset('homo.n.02')
```

```
In [36]: 강아지.lowest_common_hypernyms(인간)
```

```
Out[36]: [Synset('placental.n.01')]
```

```
In [37]: 강아지.lowest_common_hypernyms(고양이)
```

```
Out[37]: [Synset('carnivore.n.01')]
```

```
In [38]: 고양이.lowest_common_hypernyms(사자)
```

```
Out[38]: [Synset('feline.n.01')]
```

```
In [ ]:
```

육식동물

강아지

고양이과

사자

고양이

```
In [39]: 인간.min_depth()
```

```
Out[39]: 13
```

```
In [40]: 강아지.min_depth()
```

```
Out[40]: 8
```

```
In [41]: 인간.path_similarity(강아지)
```

```
Out[41]: 0.14285714285714285 | 14%
```

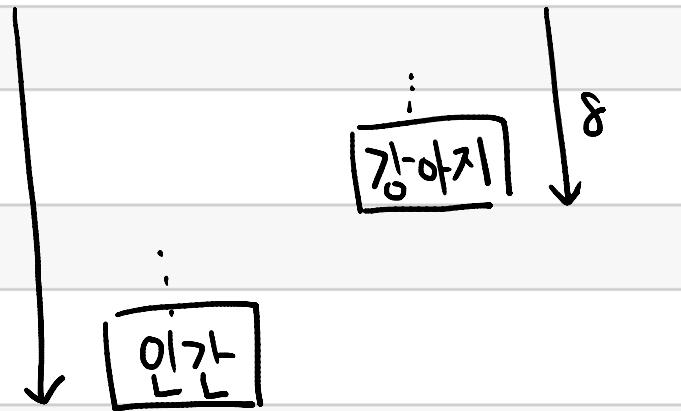
```
In [42]: 인간.path_similarity(고양이)
```

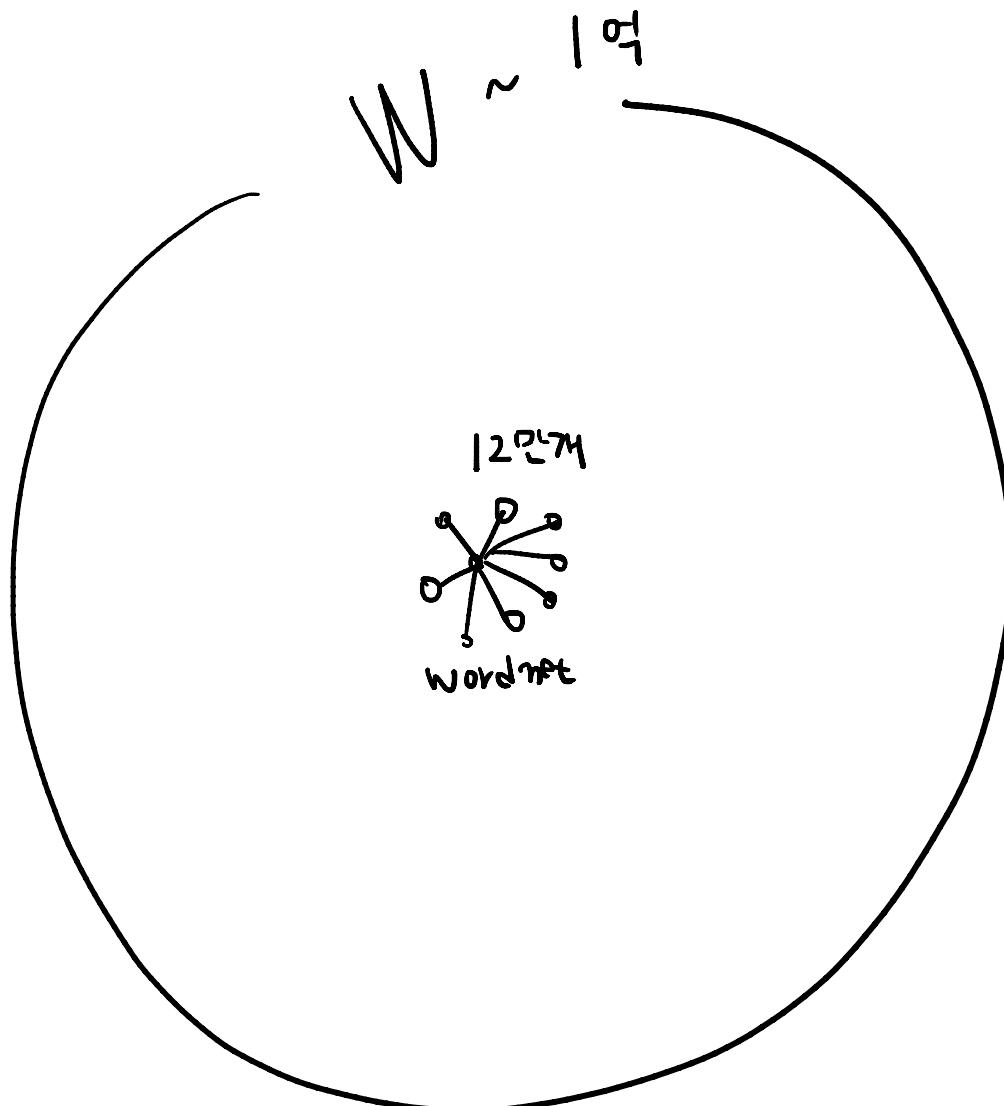
```
Out[42]: 0.14285714285714285 | 14%
```

```
In [43]: 인간.path_similarity(사자)
```

```
Out[43]: 0.125 | 12%
```

```
In [ ]:
```





정규표현식

In [45]: text = """연락주세요 010-1234-5678. 혹시 안 되면,
010.2345.7890으로 부탁드립니다."""

STD ✓ ✓
In [44]: import re Regular Expression 010-1234-5678

In [46]: 전화번호패턴 = r'\Wd{3}-\Wd{4}-\Wd{4}' "숫자3개"- "숫자4개"- "숫자4개"

In [47]: re.findall(전화번호패턴, text)

Out [47]: ['010-1234-5678']

↓
r'\Wd{3}-\Wd{4}-\Wd{4}'
num ↲

In []:

정규표현식

```
In [45]: text = """연락주세요 010-1234-5678. 혹시 안 되면,  
010.2345.7890으로 부탁드립니다."""
```

```
In [44]: import re
```

```
In [48]: 전화번호패턴 = r'\Wd{3}[.-]\Wd{4}[.-]\Wd{4}'
```

[. -] "또는"

```
In [49]: re.findall(전화번호패턴, text)
```

```
Out[49]: ['010-1234-5678', '010.2345.7890']
```

```
In [ ]:
```

Str ... ed

re.match(패턴, Str)

for text in [Str, Str, ...]:

re.match(패턴, Str)

Series .str.contains(패턴)

Str
Str
:

True

False

:

"액팅함수"



In [56]: 패턴선택 = lambda 패턴: 영단어목록[영단어목록.str.contains(패턴)]

In [57]: 패턴선택('ed\$').sample(10)

def 패턴선택(패턴):

Out[57]:

211647	unbeaued
224795	unwailed
186601	spiculated
211326	unarticled
218226	uninfuriated
12263	aralkylated
109978	magged
30437	carded
212767	uncollaged
219662	unmounted

dtype: object

return ...

"여덟 글자이면서, 세 번째 글자가 j, 여섯 번째 글자가 t인 경우"

In []: 패턴 = '^..j..t..\$'

In []:

정규표현식 META 문자

\$ " ...로 끝나는" $[^ab]$ a 또는 b가 아닌 것

^ " ...로 시작" $\{S, E\}$ 개수 지정

.

"아무 글자"

$ad | cd$ ab 또는 cd
 $()$ 그룹

[] 글자범위 ; 범위표현 e.g. 0-9, a-z, 가-힝

+

1개 이상

*

0개 이상

File Edit View Insert Cell Kernel Widgets Help

```
3    aai  
4    aali i  
dtype: object
```

In [55]: 영단어목록[영단어목록.str.contains('ed\$')].sample(10)

Out [55]:

```
212726      uncoddled  
223724      unsurnamed  
141201      peeled  
27888       busybodied  
40429       componented  
220589      unexplained  
219531      unmissed  
169041      revulsed  
218693     unjointed  
219346      unmedaled  
dtype: object
```

정규표현식 META 문자

\$ “…로 끝나는 패턴”

In []:

File Edit View Insert Cell Kernel Widgets Help

dtype: object

"여덟 글자이면서, 세 번째 글자가 j, 여섯 번째 글자가 t인 경우"

In [61]: 패턴선택('^...j...t...\$').sample(10)

Out[61]:

110321	majestic
94828	injector
2620	adjuster
50274	dejectly
129006	objectee
166286	rejecter
166291	rejector
276	abjectly
218683	unjilted
218736	unjustly

dtype: object

^ . . j . . t . . \$

In []:

File Edit View Insert Cell Kernel Widgets Help

dtype: object

"여덟 글자이면서, 세 번째 글자가 j, 여섯 번째 글자가 t인 경우"

In [63]: 패턴선택('^.j..t..\$').sample(10)

...

In [62]: 패턴선택('..j..t..').sample(10)

Out [62]: 210848 unadjectived
218675 unjesting
163758 readjuster
129029 objectization
156693 projectedly
223544 unsubjected
129016 objectioner
129004 objectionation
41385 conjecturably
2619 adjustment
dtype: object

In []:

File Edit View Insert Cell Kernel Widgets Help

"여덟 글사이면서, 세 번째 글자가 j, 여섯 번째 글자가 t인 경우"

In [63]: 패턴선택('^.j..t..\$').sample(10)

...

In [62]: 패턴선택('..j..t..').sample(10)

...

In [74]: 패턴선택('^[ghi][mno][j|k][def]\$')

^ [ghi][mno] ...

Out [74]:

78596	gold
78655	golf
86476	hold
86492	hole

dtype: object

^
\$
[ab] a 또는 b

In []:

File Edit View Insert Cell Kernel Widgets Help

In [84]: chat_words.sample(5)

Out [84]:

```
39677      dd
6735      beautiful
557       ca
8615      removing
8903      U58
dtype: object
```

"1개 이상"
m+ m
m m
m m m ...
^m+i+n+e+\$

In [85]: 패턴선택(chat_words, '^m+i+n+e+\$')

Out [85]:

```
3222          mine
29128         miiiiiiinnnnnnnnneeeeeeee
29198         miiiiiiiiiiiiinnnnnnnnneeeeeeee
29204         mmmmmmmmmiiiiiiiiinnnnnnnneeeeeeee
dtype: object
```

heellooo ...

In []:

```
In [89]: 패턴선택(chat_words, '^m*i*n*e*$')
```

만약 “도개 악성”

```
Out[89]: 36          m
          45          me
          69          i
          275         in
          1867        min
          2576        mmm
          3222        mine
          3315         n
          6946         mm
          15069        mmmm
          18067        mmmmm
          19375
          29128        miiiiiiinnnnnnnnneeeeee
          29198        miiiiiiiiiiiiinnnnnnnnneeeeee
          29204        mmmmmmmmmiiiiiiiiinnnnnnnnneeeeee
          30010          mmmmmmmmmmmmmmm
          30586          e
          31203          ne
          32129          meeeeeeeeeeee
          32158          mmmmmmmmmmmmmmm
          34375          mi
          36378          mmmmmmmmmmm
          36982          mmmmmmm
```

File Edit View Insert Cell Kernel Widgets Help

In [89]: 패턴선택(chat_words, '^m*i*n*e*\$')

...

In [91]: 패턴선택(chat_words, '^[^aeiouAEIOU]+\$').sample(10)

Out[91]: 9250 =)
20654]]]]]]]]]]]]]]]]]]
4305 ++
13471))))))))))))
10484 wz
37502 ;-(
6132 <<<<
13121 (((((
15016
43105 f|
dtype: object

NT⁵ ↘
^ [^aeiouAEIOU] + \$
부정

In []:

File Edit View Insert Cell Kernel Widgets Help

In [93]: 신문단어 = Series(t for t in 신문.words()).drop_duplicates()

In [94]: 신문단어.sample(10)

...

In [95]: 패턴선택(신문단어, '^[0123456789]+.[0123456789]+\$')

Out [95]:

159	1956
346	1953
348	1955
350	9.8
668	1997
683	160
869	400
888	8.45
891	8.47
1075	8.04
1078	7.90
1106	1.5
1116	352.7
1167	9.37
1176	9.45
1212	8.12
1217	8.14
1228	8.19

목표: 3.14 와 같은 소수정을 포함한 숫자

^ [0-9]+.\ [0-9]+\$
 ↑

In [96]: 패턴선택(신문단어, ' $^{\wedge}[0-9]+\w{.}[0-9]+\$$ ')

Out [96]:

350	9.8
888	8.45
891	8.47
1075	8.04
1078	7.90
1106	1.5
1116	352.7
1167	9.37
1176	9.45
1212	8.12
1217	8.14
1228	8.19
1231	8.22
1242	8.53
1245	8.56
1263	83.4
1487	2.80
1492	2.87
1527	3.1
2167	5.29
2173	0.7
2187	5.39

In [97]: 패턴선택(신문단어, '^[0-9]{4}\$')

Out [97]:	
159	1956
346	1953
348	1955
668	1997
2208	1986
2325	1990
2559	1989
3627	1988
3716	1985
4767	1991
6225	1977
6420	1975
6445	1981
6786	1987
8315	1929
8791	1992
9206	1917
9211	1934
12226	1979
12644	1901
12704	1972
12741	1999

$^{\text{[0-9][0-9][0-9][0-9]}}$ \$

$\text{[0-9]} *$ "0개 이상"

$+$ "1개 이상"

{n} "n개"

In [97]: 패턴선택(신문단어, '^[0-9]{4}\$')

In [98]: 패턴선택(신문단어, '^[0-9]+-[a-z]{3,5}\$')

Out [98]:

1221	30-day
1919	10-lap
7290	52-week
17616	50-state
18521	14-hour
20851	69-point
27263	10-year
34158	30-year
34191	36-day
39975	150-point
39986	20-point
40135	12-point
40214	30-point
42909	500-stock
46378	10-day
47858	90-day
47888	15-day
48892	100-share
49702	94-month

$\wedge [0-9]+-[a-z]^{3,5}\$$

{n} "n개"

{n,m} "n ≤ m"

In [97]: 패턴선택(신문단어, '^[0-9]{4}\$')

...

In [98]: 패턴선택(신문단어, '^[0-9]+-[a-z]{3,5}\$')

...

In [100]: 패턴선택(신문단어, '^[a-z]{5,}-[a-z]{2,3}-[a-z]{1,6}\$')

Out [100]:

12948	black-and-white
39511	savings-and-loan
44106	father-in-law
45462	machine-gun-totting
94123	bread-and-butter

dtype: object

{n} n개
{S,C} S≤ ≤e
{S,I} S개 이상
{,E} E개 이하

In []:

In [101]: 패턴선택(신문단어, '(ed|ing)\$')

```
C:\ProgramData\Anaconda3\lib\site-packages\ipykernel_launcher.py:  
1: UserWarning: This pattern has match groups. To actually get the  
groups, use str.extract.  
    """Entry point for launching an IPython kernel.
```

Out [101]:

28	publishing
41	Consolidated
47	named
63	used
72	caused
84	exposed
95	reported
120	causing
141	York-based
150	stopped
151	using
191	bring
212	talking
221	having
243	studied
286	led
324	making
334	replaced
372	worked

"ed, ing 끝나는"

[ed ing]

(ed | ing)\$

In [103]: 패턴선택(신문단어, '(ed|ing)\$')

ed\$ "또는" ing\$

In [104]: 패턴선택(신문단어, 'ed|ing\$')

ed "또는" ing\$

Out [104]:

28	publishing
41	Consolidated
47	named
63	used
72	caused
84	exposed
95	reported
120	causing
141	York-based
150	stopped
151	using
184	Medicine
191	bring
212	talking
221	having
243	studied
286	led
298	medical
324	making
334	replaced

File Edit View Insert Cell Kernel Widgets Help

```
45462    machine-gun-totting  
94123    bread-and-butter  
dtype: object
```

In [103]: 패턴선택(신문단어, '(ed|ing)\$')

...

In [104]: 패턴선택(신문단어, 'ed|ing\$')

...

1) 패턴확인 → 010-1234-5678

In [114]: re.findall(
r'(\w{3})-(\w{3,4})-(\w{4})',
'내 전번은 010-1234-5678입니다.')

2) () "그룹생성"

Out[114]: [('010', '1234', '5678')]

In []:

File Edit View Insert Cell Kernel Widgets Help

In [104]: 패턴선택(신문단어, 'ed|ing\$')

...

In [114]: re.findall(
r'(\Wd{3})-(\Wd{3,4})-(\Wd{4})',
'내 전번은 010-1234-5678입니다.')

"교체" substitute

Out[114]: [('010', '1234', '5678')]

In [115]: re.sub(1 2 3
r'(\Wd{3})-(\Wd{3,4})-(\Wd{4})',
r'\1-***- \3',
'내 전번은 010-1234-5678입니다.'
)

re.sub(패턴, 대체, 대상)

\1 "1번 그룹"

Out[115]: '내 전번은 010-***-5678입니다.'

In []:

File

Edit

View

Insert

Cell

Kernel

Widgets

Help

ASCII

A-Za-zA-Z

UNICODE

한글 범위

ㄱ-ㅣ	ㅏ-ㅓ	ㅓ-ㅗ	ㅗ-ㅜ
-----	-----	-----	-----

Out[115]: '내 전번은 010-****-5678입니다.'

In [117]: text = '알晦다운_우리말 ^0^ English, 123'

In [120]: re.findall('[^ㄱ-ㅣ 가-핳]+', text)

Out[120]: ['^0^', 'English,', '123']

In [118]: re.sub('[^ㄱ-ㅣ 가-핳]+', '', text)

Out[118]: '알晦다운우리말'

In [119]: re.sub('[^ㄱ-ㅣ 가-핳]+', '', text)

Out[119]: '알晦다운 우리말 '

In []:

File Edit View Insert Cell Kernel Widgets Help

Out[120]: ['^0^', 'English,', '123']

In [118]: re.sub('^[ㄱ-ㅣ 가-힣]+', '', text)

Out[118]: '알晦다운우리말'

In [119]: re.sub('^[ㄱ-ㅣ 가-힣]+', '', text)

Out[119]: '알晦다운 우리말 '

In [124]: re.sub('[ㄱ-ㅣ]+', '', 'ㅋㅋㅋ, ㅎㅎㅎ 이런 것들을 제거합니다.')

Out[124]: ' 이런 것들을 제거합니다.'

In [125]: tokens = Series([
 'processing',
 'processed',
 'processes',
 'process'
])

In [128]: tokens.str.findall('.+ing|ed|es\$')

Out[128]: 0 [(process, ing)]
 1 [(process, ed)]

File Edit View Insert Cell Kernel Widgets Help

Out[124]: ', 이런 것들을 제거합니다.'

In [125]: tokens = Series([
 'processing',
 'processed',
 'processes',
 'process'
])

(.+)(ing|ed|es)\$

do did done

In [128]: tokens.str.findall('(.+)(ing|ed|es)\$')

Out[128]: 0 [(process, ing)]
1 [(process, ed)]
2 [(process, es)]
3 []
dtype: object

타다

타고

타는

탔다 타 + ㅆ다

In []:

```
In [131]: 영화리뷰 = pd.read_csv('data/nsmc/ratings.txt', sep='\\t')
```

```
In [132]: 영화리뷰['document'].sample(10)
```

```
Out[132]: 116340
```

ㅋㅋㅋ

169853

ㅋㅋ

87978

예쁨

90039

26085 금율 삼시세끼 꼭 챙겨봅니다. 첨엔 일단 나pd님. 이서진씨 믿고본거긴한데.. 크 ...

113712 피터잭슨이 이런 000기영화를 만들다니 ㅋㅋㅋ 러블리본즈에 이어 또실망..

190524 여명은 한국에 눈이 멀었던 게지..

126063 안타까움

34649 평점 너무 약하다....견자단 완전 멋짐!!!

11666 단순한 킬링타임용 호러물이 절대 아닙니다.. 완성도도 그럴지만 무엇보다 에밀리 로즈...

Name: document, dtype: object

1) 정글화

2) 도큰화

처음에는

줄거리에 내용을 써놨네

여자 존나 짜증나네

티아 레오니가 너무

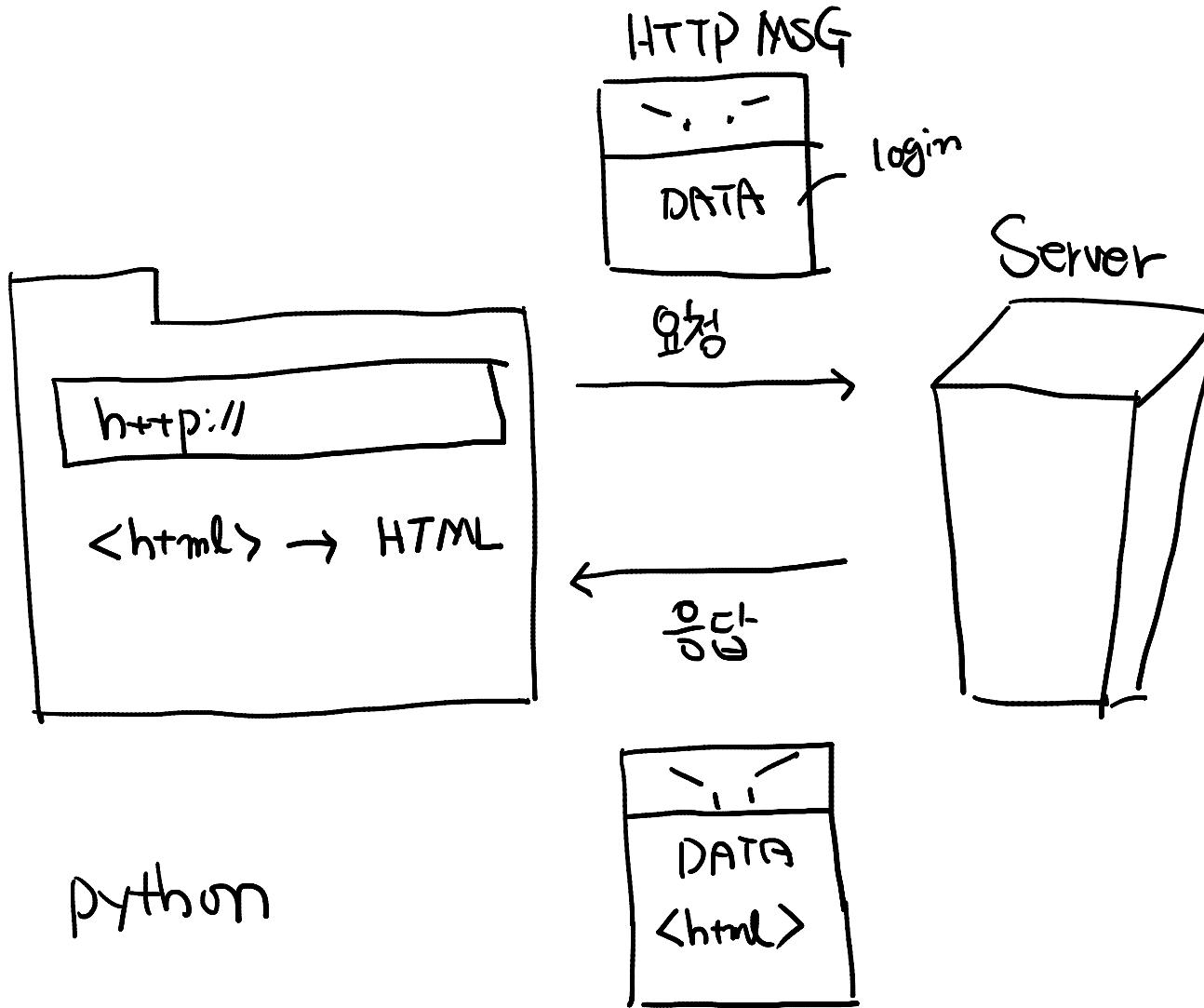
너무 보고싶어요...

첨엔 일단 나pd님.

여명은 한국에 눈이 멀었던

안타까움

평점 너무 약하다....견자단 완전



File Edit View Insert Cell Kernel Widgets Help

In [133]: `import requests`

요청

In [134]: `응답 = requests.get('http://www.example.com')`

응답

In [135]: `data = 응답.text`

응답.text

In [136]: `print(data)`



```
<!doctype html>
<html>
<head>
    <title>Example Domain</title>

    <meta charset="utf-8" />
    <meta http-equiv="Content-type" content="text/html; charset=utf-8" />
    <meta name="viewport" content="width=device-width, initial-scale=1" />
    <style type="text/css">
        body {
            background-color: #f0f0f2;
            margin: 0;
            padding: 0;
            font-family: "Open Sans", "Helvetica Neue", Helvetica, Arial, sans-serif;
```

File

Edit

View

Insert

Cell

Kernel

Widgets

Help

In [149]: html = BeautifulSoup(data, 'lxml').get(속성).text

Tag

In [150]: html.title

Out[150]: <title>NAVER</title>

Click ME!

In [151]: html.title.text

html.a

Out[151]: 'NAVER'

.find-all('a')

In [154]: links = html.find_all('a')

In [155]: for link in links:
print(link.text, link.get('href'))

상대 | 절대

뉴스스탠드 바로가기 #news_cast

주제별캐스트 바로가기 #themecast

타임스퀘어 바로가기 #time_square

쇼핑캐스트 바로가기 #shop_cast

로그인 바로가기 #account

네이버 /

네이버를 시작페이지로 http://help.naver.com/support/alias/content
s2/naverhome/naverhome_1.naver

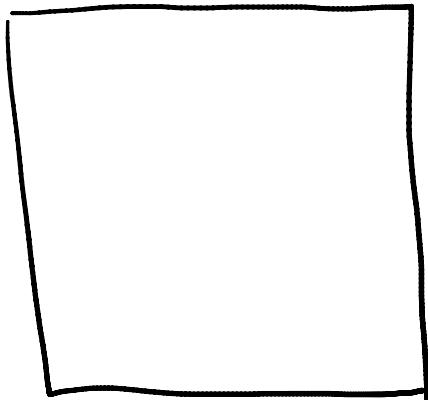
쥬니어네이버 http://jr.naver.com

해피빈 http://happybean.naver.com/main/SectionMain.nhn

http://

News

Website



기사 링크 추출

http:// ...

http:// ...



기사 본문 추출



plamtextCorpusReader
(word_tokenizer)

.words()

.sents()

File Edit View Insert Cell Kernel Widgets Help

In [179]: 응답 = requests.get(url)

In [181]: 응답.encoding = 'utf-8'

In [183]: 웹페이지 = BeautifulSoup(응답.text, 'lxml')

In [185]: 웹페이지.title.text

Out[185]: '방남 北응원단 "기존에 없던 것 보여줄 생각...보시면 압네다"(종합)'

In [188]: 본문 = '\n'.join(
 문단.text for 문단 in 웹페이지.select('.article p'))

CSS 선택자

<div class="article">

In [189]: print(본문)

<p> ...</p>

(도라산·서울=연합뉴스) 공동취재단 이정진 김정은 기자 = 평창동 계올림픽에서 남북 선수들을 응원할 북한 응원단과 태권도시범단 등이 7일 방남했다.

버스 9대에 나눠타고 이날 오전 9시 28분 경기 파주의 도라산 남북 출입사무소(CIQ)에 도착한 이들은 10시 13분부터 기자단, 응원단, 태권도시범단, 민족올림픽위원회(NOC) 관계자 순으로 차례차례 남측 출구로 나왔다.

남성들은 검은색 코트에 털모자 여성들은 붉은 코트에 검은색 털

File Edit View Insert Cell Kernel Widgets Help

```
In [200]: with open(filename, 'w', encoding='utf-8') as file:  
    file.write(내용)
```

```
In [201]: from nltk.corpus import PlaintextCorpusReader
```

```
In [202]: processor = pyko.OpenKoreanTextProcessor()
```

```
In [203]: reader = PlaintextCorpusReader(  
            root='.',  
            fileids=['연합뉴스-0505000000AKR20180207084551014.txt'],  
            word_tokenizer=processor  
)
```

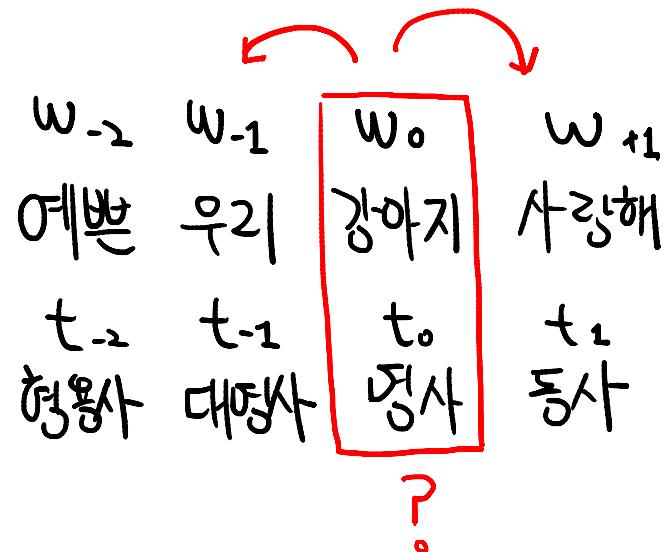
```
In [204]: tokens = reader.words()
```

```
In [205]: tokens
```

```
Out[205]: [[', '방남', '北', '응원단', "", '기존', '에', '없던', '것',  
...]]
```

```
In [206]: nltk.Text(tokens).collocations()
```

모찌



현실:
→ 내째무 거 보 려 보

In [212]: `sejong.sents()`

$[[w_1, w_2, w_3], [], []]$

Out [212]: `[['뭐', '타고', '가?'], ['지하철.'], ['기차?'], ['아침에', '몇', '시에', '타고', '가는데?'], ['아침에', '한', '일찍', '가면은', '일곱시', '이십분'], ...], ['음.'], ...]`

가? 가/VV+ㅏ/EF+?/SF

In [210]: `tagged_sents = sejong.sents(tagged=True)`

H./T,+H2/T2+H3/T3

In [211]: `tagged_sents`

$[[(w_1, T_1), (w_2, T_2), \dots], \dots]$

Out [211]: `[[('뭐', 'NP'), ('ㅌ', 'VV'), ('고', 'EC'), ('가', 'VV'), ('ㅏ', 'EF'), ('?', 'SF')], [('지하철', 'NNG'), ('.', 'SF')], [('기차', 'NNG'), ('?', 'SF')], [('아침', 'NNG'), ('에', 'JKB'), ('몇', 'M'), ('시', 'NNB'), ('에', 'JKB'), ('ㅌ', 'VV'), ...], [('아침', 'NNG'), ('에', 'JKB'), ('한', 'MM'), ('일찍', 'MAG'), ('가', 'V'), ('면은', 'EC'), ...], [('음', 'IC'), ('.', 'SF')], ...]`In [214]: `pd.read_csv('data/sejong_tagset.txt', index_col='Tag')`

Out [214]:

Description

Tag

NNG

일반 명사

BNC

Sejong

<w pos="NOUN" >dog</w> <w>

• tagged-words()

→ (dog, NOUN)

강아지 강아지/NNG

• words(tagged=True)

→ (강아지, NNG)

품사 분류 (tagset)

영어

Class	POS
Open	Noun
	Verb
	Adjective
	Adverb
	Interjection
Close	Pronoun
	Preposition
	Conjunction
	Article

한국어

5언	9품사
체언	명사
	대명사
	수사
용언	동사
	형용사
수식언	관형사
	부사
관계언	조사
독립언	감탄사

File Edit View Insert Cell Kernel Widgets Help

In [214]: pd.read_csv('data/sejong_tagset.txt', index_col='Tag')

Out [214]:

Tag	Description
NNG	일반 명사
NNP	고유 명사
NNB	의존 명사
NR	수사
NP	대명사
VV	동사
VA	형용사
VX	보조 용언
VCP	긍정 지정사
VCN	부정 지정사
MM	관형사
MAG	일반 부사

File Edit View Insert Cell Kernel Widgets Help

곱시', '이십분', ...], ['음.'], ...]

In [210]: tagged_sents = sejong.sents(tagged=True)

In [211]: tagged_sents → [NP, VV, EC, VV, EF, ...]

Out[211]: [[('뭐', 'NP'), ('타', 'VV'), ('고', 'EC'), ('가', 'VV'), ('ㅏ', 'EF'), ('?', 'SF')], [('지하철', 'NNG'), ('.', 'SF')], [('기차', 'NNG'), ('?', 'SF')], [('아침', 'NNG'), ('에', 'JKB'), ('몇', 'M'), ('시', 'NNB'), ('에', 'JKB'), ('타', 'VV'), ...], [('아침', 'NNG'), ('에', 'JKB'), ('한', 'MM'), ('일찍', 'MAG'), ('가', 'V'), ('면은', 'EC'), ...], [('음', 'IC'), ('.', 'SF')], ...]

In [214]: pd.read_csv('data/sejong_tagset.txt', index_col='Tag')

...

In []:

```
In [215]: tags = []
for sent in tagged_sents:
    for token in sent:
        tag = token[1]
        tags.append(tag)
```

tags 만들기

```
In [219]: tag_pairs = list(nltk.bigrams(tags))
```

```
In [220]: tag_pairs[:10]
```

```
Out[220]: [('NP', 'VV'),
            ('VV', 'EC'),
            ('EC', 'VV'),
            ('VV', 'EF'),
            ('EF', 'SF'),
            ('SF', 'NNG'),
            ('NNG', 'SF'),
            ('SF', 'NNG'),
            ('NNG', 'SF'),
            ('SF', 'NNG')]
```

NP, VV, EC

(NP, VV), (VV, EC)

(?, NN.)

```
In [ ]:
```

```
In [220]: tag_pairs[:10]
```

...

```
In [221]: target_tags = []
for t1, t2 in tag_pairs:
    if t2.startswith('NN'):
        target_tags.append(t1)
```

^ NN

```
In [224]: Series(target_tags).value_counts()[:10]
```

Out[224]:

ETM	53694
SP	34334
MM	34006
NNG	32204
SF	22829
MAG	21477
IC	11725
EC	10918
NR	9261
JKB	8964

dtype: int64

수집

→ SP
→ NNG

나, 이, 우리

우리 강아지

```
In [ ]:
```

In [227]: words, tags = zip(*[token for sent in tagged_sents for token in sent])

In [228]: tags = Series(tags, index=words)

In [229]: tags[:10]

Out [229]: 워 NP
타 VV
고 EC
가 VV
ㅏ EF
? SF
지하철 NNG
. SF
기차 NNG
? SF
dtype: object

In []:

(워, NP), (타, VV),

[워, 타, ...]
[NP, VV, ...]

File Edit View Insert Cell Kernel Widgets Help

In [237]: 기본태그기.tag(sample)

Out[237]: [('뭐', 'NNG'),
 ('타', 'NNG'),
 ('고', 'NNG'),
 ('가', 'NNG'),
 ('ㅏ', 'NNG'),
 ('?', 'NNG')]

In [239]: tagged_sents[0]

Out[239]: [('뭐', 'NP'), ('타', 'VV'), ('고', 'EC'), ('가', 'W'), ('ㅏ', 'E'), ('?', 'SF')]

In [238]: 기본태그기.evaluate(tagged_sents[:1000])

Out[238]: 0.11324162845630123

In []:

Out[240]: 0.11524102040000120

조건부 빈도집계

```
In [240]: 단어별품사빈도 = nltk.ConditionalFreqDist(  
          token for sent in tagged_sents for token in sent  
         )
```

```
In [241]: 단어별품사빈도
```

```
Out[241]: ConditionalFreqDist(nltk.probability.FreqDist,  
                           {'뭐': FreqDist({'IC': 7321,  
                                         'MAG': 1,  
                                         'NNG': 11,  
                                         'NNP': 1,  
                                         'NP': 5284,  
                                         'UNT': 1}),  
                            '타': FreqDist({'MM': 3,  
                                         'NNB': 21,  
                                         'NNG': 11,  
                                         'UNT': 3,  
                                         'VV': 570,  
                                         'XPN': 1}),  
                            '고': FreqDist({'EC': 19104,  
                                         'EF': 518,  
                                         'ETM': 1,  
                                         'IC': 9,  
                                         'JC': 2,  
                                         'JKC': 12})}
```

특정단어 가장높은 도수 품사

In [249]: 품사추출 = dict(
↓
(word, 단어별품사빈도[word].max()) for word in 단어도수.index)

LookUp

In [251]: 품사추출
(, , 'SP')

Out [251]: {':': 'SP',
'.': 'SF',
'하': 'VV',
'이': 'JKS',
'는': 'ETM',
'가': 'JKS',
'?': 'SF',
'에': 'JKB',
'거': 'NNB',
'어': 'IC',
'ㄴ': 'ETM',
'고': 'EC',
'을': 'JKO',
'있': 'VA',
'(으)': 'VCP',
'그': 'MM',
'은': 'JX',
'되': 'VV',
'아': 'IC'}

File Edit View Insert Cell Kernel Widgets Help

In [247]: 단어도수 = Series(t[0] for t in 토큰생성기()).value_counts()

In [248]: 단어도수[:10]

...

In [249]: 품사추측 = dict(
(word, 단어별품사빈도[word].max()) for word in 단어도수.index)

In [251]: 품사추측

기계학습모델

...



In [252]: baseline_tagger = nltk.UnigramTagger(model=품사추측)

In [*]: baseline_tagger.evaluate(tagged_sents)

In []:

File Edit View Insert Cell Kernel Widgets Help

...

In [252]: `baseline_tagger = nltk.UnigramTagger(model=품사추측)`

In [253]: `baseline_tagger.evaluate(tagged_sents)`

Out [253]: 0.8603526460812297 ← 정확도

In [254]: `sample = [token[0] for token in tagged_sents[3]]`
sample

설정 예측 ...
↓ ↓

In [255]: `baseline_tagger.tag(sample)`

...

In [256]: `tagged_sents[3]`

...

File Edit View Insert Cell Kernel Widgets Help

Out [253]: 0.8603526460812297

In [254]: sample = [token[0] for token in tagged_sents[3]]
sample

...

N=1

1-gram unigram

In [255]: baseline_tagger.tag(sample)

...

2-gram bigram

In [256]: tagged_sents[3]

...

3-gram

N-1

w₋₂ w₋₁ w₀ w₁

t₋₂ t₋₁ t₀ w₁

N-Gram Tagging

In []:

```
In [255]: baseline_tagger.tag(sample)
```

```
In [256]: tagged_sents[3]
```

train

test

N-Gram Tagging

```
In [257]: unigram_tagger = nltk.UnigramTagger(tagged_sents)
```

```
In [*]: unigram_tagger.evaluate(tagged_sents) "평가"
```

```
In [ ]:
```

"전" ~ "한"



```
In [265]: bigram_tagger.evaluate(test_sents)
```

```
Out[265]: 0.6856353591160221
```

t_{-1} t_0

```
In [266]: t0 = nltk.DefaultTagger('NNG')
```

```
t1 = nltk.UnigramTagger(train_sents, backoff=t0)
```

```
t2 = nltk.BigramTagger(train_sents, backoff=t1)
```

None

None

```
In [267]: t2.evaluate(test_sents)
```

```
Out[267]: 0.9252839165131983
```

```
In [268]: t2.tag(sample)
```

```
Out[268]: [('아침', 'NNG'),  
           ('에', 'JKB'),  
           ('몇', 'MM'),  
           ('시', 'NNB'),  
           ('에', 'JKB'),  
           ('타', 'VV'),  
           ('고', 'EC'),  
           ('가', 'VV'),  
           ('는데', 'EC'),  
           ('?', 'SF')]
```

File Edit View Insert Cell Kernel Widgets Help

In [270]: `import pickle`

In [271]: `with open('sejong_tagger.pkl', 'wb') as file:`
 `pickle.dump(t2, file)`

In [272]: `del t2`

In [274]: `with open('sejong_tagger.pkl', 'rb') as file:`
 `t2 = pickle.load(file)`

In [275]: `t2.tag(sample)`

Out [275]: `[('아침', 'NNG'),
 ('에', 'JKB'),
 ('몇', 'MM'),
 ('시', 'NNB'),
 ('에', 'JKB'),
 ('타', 'VV'),
 ('고', 'EC'),
 ('가', 'VV'),
 ('는데', 'EC'),
 ('?', 'SF')]`

