

▼ 공백 기반 토큰화

한 문장에서 단어의 수는 어떻게 정의할 수 있을까요? "그녀는 나와 밥을 먹는다" 라는 문장이 주어지면 공백 기준으로 나누어 1: 그녀는 2: 나와 3: 밥을 4: 먹는다 4개 단어로 이루어졌다고 단정 지을 수 있을까요? 어쩌면 1: 그녀 2: 는 3: 나 4: 와 5: 밥 6: 을 7: 먹는다 처럼 잘게 잘게 쪼개어 7개 단어로 이루어졌다고 할 수도 있겠죠? 그것은 우리가 정의할 **토큰화 기법이 결정할 부분**입니다!

문장을 어떤 기준으로 쪼개었을 때, 쪼개진 각 단어들을 **토큰(Token)** 이라고 부릅니다. 그리고 그 쪼개진 기준이 **토큰화(Tokenization) 기법**에 의해 정해지죠. 이번 스텝에서는 토큰화의 여러 가지 기법에 대해 배워보도록 하겠습니다.

자연어의 노이즈를 제거하는 방법 중 하나로 우리는 `Hi`, `를 Hi` 와 `,` 로 나누기 위해 문장부호 양옆에 공백을 추가해 주었습니다. 그것은 이 **공백 기반 토큰화**를 사용하기 위해서였죠! 당시의 예제 코드를 다시 가져와 공백을 기반으로 토큰화를 진행해 보겠습니다.

```
1 corpus = \
2 """
3 in the days that followed i learned to spell in this uncomprehending way a great
4 but my teacher had been with me several weeks before i understood that everythin
5 one day , we walked down the path to the well house , attracted by the fragran
6 some one was drawing water and my teacher placed my hand under the spout .
7 as the cool stream gushed over one hand she spelled into the other the word wate
8 i stood still , my whole attention fixed upon the motions of her fingers .
9 suddenly i felt a misty consciousness as of something forgotten a thrill of retu
10 i knew then that w a t e r meant the wonderful cool something that was flowing
11 that living word awakened my soul , gave it light , hope , joy , set it free
12 there were barriers still , it is true , but barriers that could in time be sw
13 """
14 # HINT : split()을 사용하여 공백토큰화를 수행하세요.
15 tokens = corpus.split()
16
17 print("문장이 포함하는 Tokens:", tokens)
```



문장이 포함하는 Tokens: ['in', 'the', 'days', 'that', 'followed', 'i', 'learned',

▼ 형태소 기반 토큰화

하지만 우리에게 영어 문장이 아닌 한국어 문장을 처리할 일이 더 많을 것이고, 한국어 문장은 **공백 기준**으로 토큰화를 했다간 **영망진창의 단어들이 등장하는 것**을 알 수 있습니다. 문장부호처럼 "**은 / 는 / 이 / 가**" 양옆에 공백을 붙이자구요? 글썄요... 가로 시작하는 단어만 해도 가면, 가위, 가족, 가수... 의도치 않은 변형이 너무나도 많이 일어날 것 같네요!


```
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.7/
Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.7/dist-p
Requirement already satisfied: urllib3!=1.25.0,!1.25.1,<1.26,>=1.21.1 in /usr
Requirement already satisfied: PySocks!=1.5.7,>=1.5.6 in /usr/local/lib/pythor
Installing collected packages: JPype1, colorama, beautifulsoup4, konlpy
  Attempting uninstall: beautifulsoup4
    Found existing installation: beautifulsoup4 4.6.3
    Uninstalling beautifulsoup4-4.6.3:
      Successfully uninstalled beautifulsoup4-4.6.3
  Successfully installed JPype1-1.3.0 beautifulsoup4-4.6.0 colorama-0.4.4 konlpy
```

```
1 cd Mecab-ko-for-Google-Colab/
```

```
/content/Mecab-ko-for-Google-Colab
```

```
1 ! bash install_mecab-ko_on_colab190912.sh
```

```
Installing konlpy.....
Requirement already satisfied: konlpy in /usr/local/lib/python3.7/dist-package
Requirement already satisfied: JPype1>=0.7.0 in /usr/local/lib/python3.7/dist-
Requirement already satisfied: lxml>=4.1.0 in /usr/local/lib/python3.7/dist-pa
Requirement already satisfied: numpy>=1.6 in /usr/local/lib/python3.7/dist-pac
Requirement already satisfied: colorama in /usr/local/lib/python3.7/dist-packe
Requirement already satisfied: beautifulsoup4==4.6.0 in /usr/local/lib/python3
Requirement already satisfied: tweepy>=3.7.0 in /usr/local/lib/python3.7/dist-
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.7/d
Requirement already satisfied: requests[socks]>=2.11.1 in /usr/local/lib/pythc
Requirement already satisfied: requests-oauthlib>=0.7.0 in /usr/local/lib/pyth
Requirement already satisfied: six>=1.10.0 in /usr/local/lib/python3.7/dist-pa
Requirement already satisfied: oauthlib>=3.0.0 in /usr/local/lib/python3.7/dis
Requirement already satisfied: urllib3!=1.25.0,!1.25.1,<1.26,>=1.21.1 in /usr
Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.7/d
Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.7/dist-p
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.7/
Requirement already satisfied: PySocks!=1.5.7,>=1.5.6 in /usr/local/lib/pythor
Done
Installing mecab-0.996-ko-0.9.2.tar.gz.....
Downloading mecab-0.996-ko-0.9.2.tar.gz.....
from https://bitbucket.org/eunjeon/mecab-ko/downloads/mecab-0.996-ko-0.9.2.tar
--2021-11-26 07:19:21-- https://bitbucket.org/eunjeon/mecab-ko/downloads/mecab-0.996-ko-0.9.2.tar
Resolving bitbucket.org (bitbucket.org)... 104.192.141.1, 2406:da00:ff00::3403
Connecting to bitbucket.org (bitbucket.org)|104.192.141.1|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://bbuseruploads.s3.amazonaws.com/eunjeon/mecab-ko/downloads/me
--2021-11-26 07:19:22-- https://bbuseruploads.s3.amazonaws.com/eunjeon/mecab-
Resolving bbuseruploads.s3.amazonaws.com (bbuseruploads.s3.amazonaws.com)... 5
Connecting to bbuseruploads.s3.amazonaws.com (bbuseruploads.s3.amazonaws.com)|
HTTP request sent, awaiting response... 200 OK
Length: 1414979 (1.3M) [application/x-tar]
Saving to: 'mecab-0.996-ko-0.9.2.tar.gz'

mecab-0.996-ko-0.9. 100%[=====>] 1.35M 2.78MB/s in 0.5s

2021-11-26 07:19:23 (2.78 MB/s) - 'mecab-0.996-ko-0.9.2.tar.gz' saved [1414979]

Done
Unpacking mecab-0.996-ko-0.9.2.tar.gz.....
Done
```

```

Change Directory to mecab-0.996-ko-0.9.2.....
installing mecab-0.996-ko-0.9.2.tar.gz.....
configure
make
make check
make install
ldconfig
Done
Change Directory to /content
Downloading mecab-ko-dic-2.1.1-20180720.tar.gz.....
from https://bitbucket.org/eunjeon/mecab-ko-dic/downloads/mecab-ko-dic-2.1.1-2021-11-26\_07:20:55--https://bitbucket.org/eunjeon/mecab-ko-dic/downloads/
Resolving bitbucket.org (bitbucket.org)... 104.192.141.1, 2406:da00:ff00::22c5
Connecting to bitbucket.org (bitbucket.org)|104.192.141.1|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://bbuseruploads.s3.amazonaws.com/a4fcd83e-34f1-454e-a6ac-c242c--2021-11-26\_07:20:55--https://bbuseruploads.s3.amazonaws.com/a4fcd83e-34f1-

```

```
1 from konlpy.tag import Mecab
```

```
1 mecab = Mecab()
```

```
1 # None자리에 문장을 넣어보고 토큰화 결과를 출력해보세요.
```

```
2
```

```
3 # 예시문장 : 자연어처리가너무재밌어서밥먹는것도가끔까먹어요
```

```
4 print(mecab.morphs('자연어처리가너무재밌어서밥먹는것도가끔까먹어요'))
```

```
['자연어', '처리', '가', '너무', '재밌', '어서', '밥', '먹', '는', '것', '도', '가끔', '까', '먹', '어', '요']
```

```
1 from konlpy.tag import Hannanum,Kkma,Komoran,Mecab,Okt
```

```
2 tokenizer_list = [Hannanum(),Kkma(),Komoran(),Mecab(),Okt()]
```

```
3
```

```
4 kor_text = '코로나바이러스는 2019년 12월 중국 우한에서 처음 발생한 뒤 전 세계로 확산된, 새로운 유형
```

```
5
```

```
6 for tokenizer in tokenizer_list:
```

```
7     print('[{}] \n{}'.format(tokenizer.__class__.__name__, tokenizer.pos(kor_text)))
```

```
[Hannanum]
```

```
[('코로나바이러스', 'N'), ('는', 'J'), ('2019년', 'N'), ('12월', 'N'), ('중국', 'N'), ('', 'N')]
```

```
[Kkma]
```

```
[('코로나', 'NNG'), ('바', 'NNG'), ('이러', 'MAG'), ('슬', 'VV'), ('는', 'ETD'), ('', 'N')]
```

```
[Komoran]
```

```
[('코로나바이러스', 'NNP'), ('는', 'JX'), ('2019', 'SN'), ('년', 'NNB'), ('12월', 'N'), ('', 'N')]
```

```
[Mecab]
```

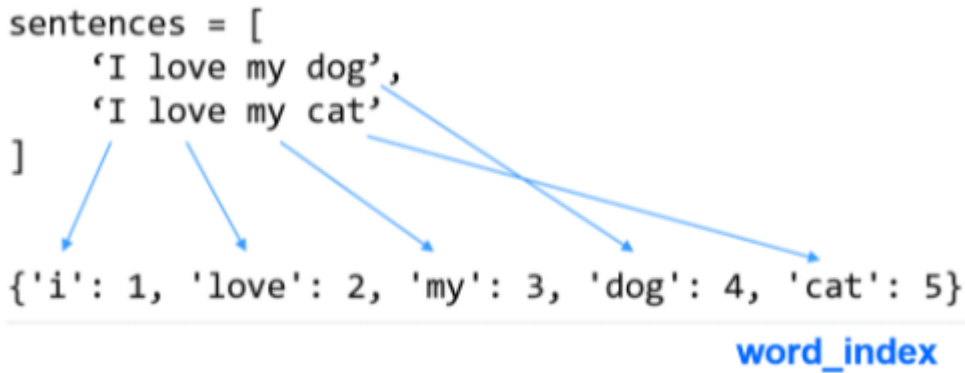
```
[('코로나', 'NNP'), ('바이러스', 'NNG'), ('는', 'JX'), ('2019', 'SN'), ('년', 'NNBC'), ('', 'N')]
```

```
[Okt]
```

```
[('코로나바이러스', 'Noun'), ('는', 'Josa'), ('2019년', 'Number'), ('12월', 'Number'), ('', 'Noun')]
```

▼ 인코딩

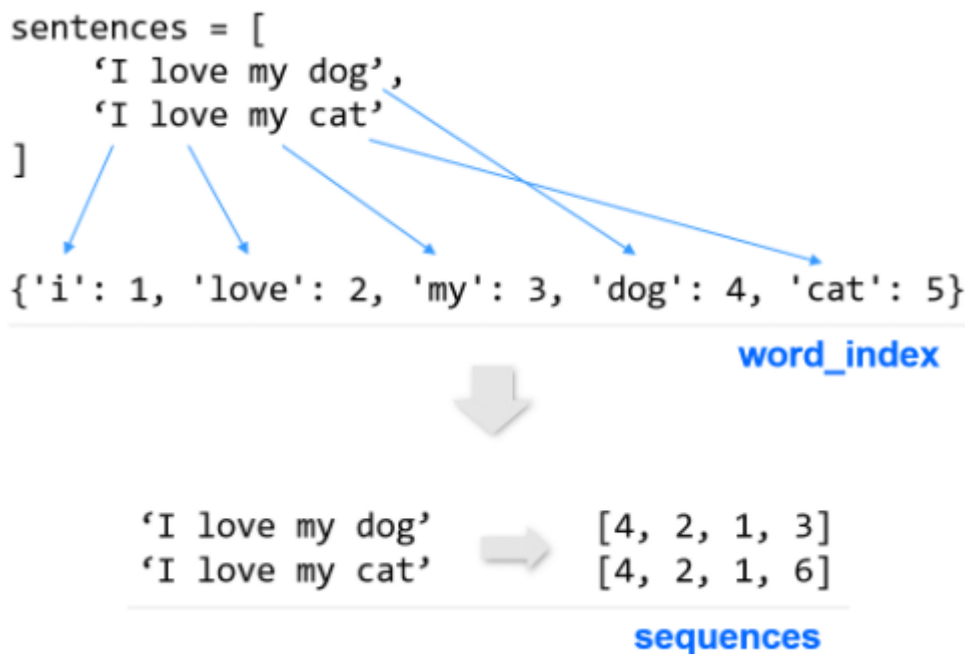
▼ 단어 기반 인코딩



```
1 import tensorflow as tf
2 from tensorflow import keras
3 from tensorflow.keras.preprocessing.text import Tokenizer
4
5 sentences = [
6     'I love my dog',
7     'I love my cat',
8 ]
9
10 tokenizer = Tokenizer(num_words = 100)
11 tokenizer.fit_on_texts(sentences) # 문자 데이터를 입력받아서 리스트의 형태로 변환
12 word_index = tokenizer.word_index # 토큰별 단어에 index를 매핑시켜준다.
13 print(word_index)
```

```
{'i': 1, 'love': 2, 'my': 3, 'dog': 4, 'cat': 5}
```

▼ 텍스트를 시퀀스로 변환하기



```
1 import tensorflow as tf
2 from tensorflow import keras
```

```

3 from tensorflow.keras.preprocessing.text import Tokenizer
4
5 sentences = [
6     'I love my dog',
7     'I love my cat',
8     'You love my dog!',
9     'Do you think my dog is amazing?'
10 ]
11
12 tokenizer = Tokenizer(num_words = 100)
13 tokenizer.fit_on_texts(sentences)
14 word_index = tokenizer.word_index
15
16 sequences = tokenizer.texts_to_sequences(sentences) # 텍스트를 시퀀스로 변환
17
18 print(word_index)
19 print(sequences)

```

{'my': 1, 'love': 2, 'dog': 3, 'i': 4, 'you': 5, 'cat': 6, 'do': 7, 'think': 8, 'amazing': 9, 'is': 10}

 [[4, 2, 1, 3], [4, 2, 1, 6], [5, 2, 1, 3], [7, 5, 8, 1, 3, 9, 10]]

▼ 패딩 설정하기

[[5, 3, 2, 4], [5, 3, 2, 7], [6, 3, 2, 4], [8, 6, 9, 2, 4, 10, 11]]

sequences



```

[[ 0  0  0  5  3  2  4]
 [ 0  0  0  5  3  2  7]
 [ 0  0  0  6  3  2  4]
 [ 8  6  9  2  4 10 11]]

```

padded

```

1 import tensorflow as tf
2 from tensorflow import keras
3 from tensorflow.keras.preprocessing.text import Tokenizer
4 from tensorflow.keras.preprocessing.sequence import pad_sequences
5
6 sentences = [
7     'I love my dog',
8     'I love my cat',
9     'You love my dog!',
10    'Do you think my dog is amazing?'
11 ]
12
13 tokenizer = Tokenizer(num_words = 100, oov_token="<OOV>")
14 tokenizer.fit_on_texts(sentences)
15 word_index = tokenizer.word_index

```

```

16 sequences = tokenizer.texts_to_sequences(sentences)
17
18 padded = pad_sequences(sequences)
19
20 print(word_index)
21 print(sequences)
22 print(padded)

{'<OOV>': 1, 'my': 2, 'love': 3, 'dog': 4, 'i': 5, 'you': 6, 'cat': 7, 'do': 8}
[[5, 3, 2, 4], [5, 3, 2, 7], [6, 3, 2, 4], [8, 6, 9, 2, 4, 10, 11]]
[[ 0  0  0  5  3  2  4]
 [ 0  0  0  5  3  2  7]
 [ 0  0  0  6  3  2  4]
 [ 8  6  9  2  4 10 11]]

```

▼ 한글-영어 코퍼스(실제 데이터로)로 토큰화수행하기

```

1 import os
2
3 import matplotlib.pyplot as plt
4 import tensorflow as tf
5 import numpy as np

1 path_to_file = '/content/drive/MyDrive/Colab Notebooks/영우4기_자연어/dataset/korean'

1 with open(path_to_file, "r", encoding='utf-8') as f:
2     raw = f.read().splitlines()
3
4 print("Data Size", len(raw))
5 print("Example:")
6 for sen in raw[0:100][::20]: print(">>", sen)

Data Size 94123
Example:
>> 개인용 컴퓨터 사용의 상당 부분은 "이것보다 뛰어날 수 있느냐?"
>> 북한의 핵무기 계획을 포기하도록 하려는 압력이 거세지고 있는 가운데, 일본과 북한의 외교관들이 외교
>> "경호 로봇가 침입자나 화재를 탐지하기 위해서 개인적으로, 그리고 전문적으로 사용되고 있습니다."
>> 수자원부 당국은 논란이 되고 있고, 막대한 비용이 드는 이 사업에 대해 내년에 건설을 시작할 계획이다.
>> 또한 근력 운동은 활발하게 걷는 것이나 최소한 20분 동안 뛰는 것과 같은 유산소 활동에서 얻는 운동

```

```

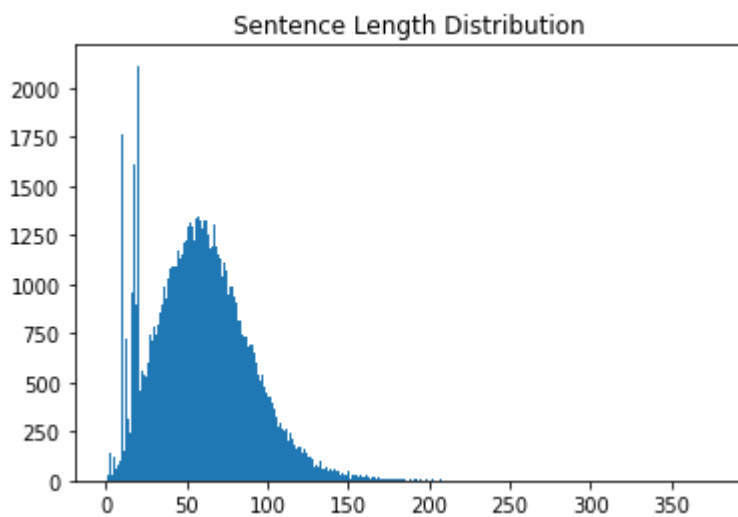
1 min_len = 999
2 max_len = 0
3 sum_len = 0
4
5 for sen in raw:
6     length = len(sen)
7     if min_len > length : min_len = length
8     if max_len < length : max_len = length
9     sum_len += length
10
11 print("문장의 최단 길이 :",min_len)

```

```
12 print("문장의 최장 길이 :", max_len)
13 print("문장의 평균 길이 :", sum_len//len(raw))
```

```
문장의 최단 길이 : 1
문장의 최장 길이 : 377
문장의 평균 길이 : 60
```

```
1 sentence_length = np.zeros((max_len), dtype=np.int)
2 for sen in raw:
3     sentence_length[len(sen)-1] += 1
4
5 plt.bar(range(max_len), sentence_length, width =1.0)
6 plt.title("Sentence Length Distribution")
7 plt.show()
```



```
1 def check_sentence_with_length(raw, length):
2     count = 0
3
4     for sen in raw:
5         if len(sen) == length:
6             print(sen)
7             count += 1
8             if count > 100: return
9
10 check_sentence_with_length(raw, 1)
```

```
,
```

```
1 for idx, _sum in enumerate(sentence_length):
2     if _sum > 1500:
3         print("Outlier Index :", idx+1)
```

```
Outlier Index : 11
Outlier Index : 19
Outlier Index : 21
```

```
1 check_sentence_with_length(raw, 11)
```



```

1 min_len = 999
2 max_len = 0
3 sum_len = 0
4
5 cleaned_corpus = list(set(raw)) # set 사용해서 중복 제거!!
6 print("Data Size :", len(cleaned_corpus))
7
8 for sen in cleaned_corpus:
9     length = len(sen)
10    if min_len > length : min_len = length
11    if max_len < length : max_len = length
12    sum_len += length
13
14 print("문장의 최단 길이 :", min_len)
15 print("문장의 최장 길이 :", max_len)
16 print("문장의 평균 길이 :", sum_len//len(cleaned_corpus))

```

```

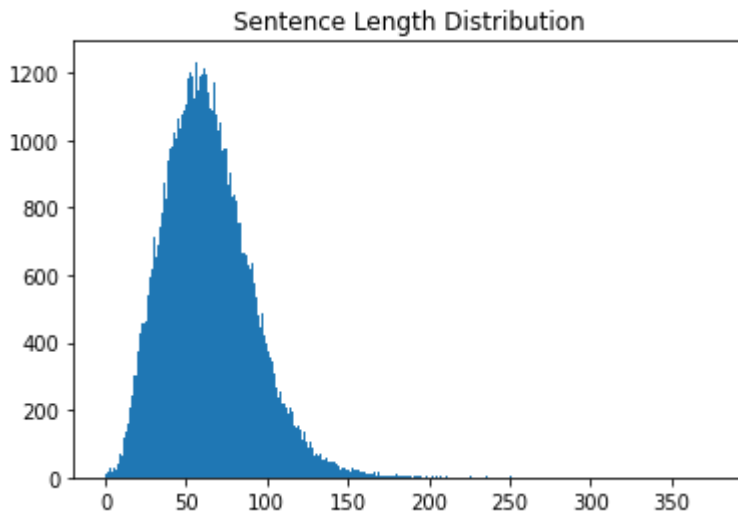
Data Size : 77591
문장의 최단 길이 : 1
문장의 최장 길이 : 377
문장의 평균 길이 : 64

```

```

1 sentence_length = np.zeros((max_len), dtype=np.int)
2 for sen in cleaned_corpus: # 중복이 제거된 코퍼스
3     sentence_length[len(sen)-1] += 1
4
5 plt.bar(range(max_len), sentence_length, width =1.0)
6 plt.title("Sentence Length Distribution")
7 plt.show()

```



```

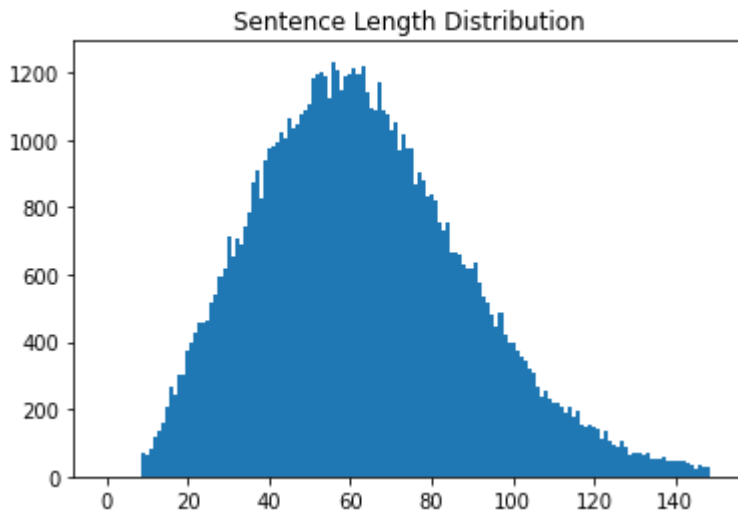
1 max_len = 150
2 min_len = 10
3
4 filtered_corpus = [s for s in cleaned_corpus if (len(s) < max_len) & (len(s) >=
5
6 sentence_length = np.zeros((max_len), dtype=np.int)
7 for sen in filtered_corpus: # 필터가 적용된 코퍼스
8     sentence_length[len(sen)-1] += 1
9

```

```

10 plt.bar(range(max_len), sentence_length, width =1.0)
11 plt.title("Sentence Length Distribution")
12 plt.show()

```



▼ 공백 기반 토큰화

```

1 # Quiz : 정제된 데이터를 공백 기반으로 토큰화하여 list에 저장한 후, tokenize()함수를 사용해 단어 .
2 # 그리고 단어 사전의 크기를 확인하세요.
3
4 def tokenize(corpus):
5
6     tokenizer = tf.keras.preprocessing.text.Tokenizer(filters='')
7     tokenizer.fit_on_texts(corpus) # 문자 -> 리스트
8
9     tensor = tokenizer.texts_to_sequences(corpus) # 텍스트 -> 시퀀스
10    tensor = tf.keras.preprocessing.sequence.pad_sequences(tensor, padding='post')
11
12    return tensor, tokenizer

```

1 # 정제된 데이터를 공백 기반으로 토큰화하여 저장하는 코드를 직접 작성해보세요.

```

2 split_corpus = []
3
4 for kor in filtered_corpus:
5     split_corpus.append(kor.split())
6     # 코드작성해주세요.

```

```

1 split_tensor, split_tokenizer = tokenize(split_corpus)
2 print("Split Vocab Size :", len(split_tokenizer.index_word))

```

Split Vocab Size : 237435

```

1 for idx, word in enumerate(split_tokenizer.word_index):
2     print(idx, ":", word)
3
4     if idx > 10: break

```

```

0 : 이
1 : 밝혔다.
2 : 있다.
3 : 말했다.
4 : 수
5 : 있는
6 : 그는
7 : 대한
8 : 위해
9 : 전했다.
10 : 지난
11 : 이번

```

▼ 형태소 기반 토큰화

```

1 # 위에서 사용한 코드를 활용해 Mecab단어 사전을 만들어주세요.
2 # Hint mecab.morphs() --> 형태소 분석 수행
3
4 def mecab_split(sentence): # 형태소 분석역할
5     return mecab.morphs(sentence)
6
7 mecab_corpus = []
8
9 # mecab 단어장 생성
10 for kor in filtered_corpus:
11     mecab_corpus.append(mecab_split(kor)) # 형태소 분석 '텍스트'형태

1 mecab_tensor, mecab_tokenizer = tokenize(mecab_corpus) # 텍스트 ---> 숫자
2 print("Mecab Vocab size :", len(mecab_tokenizer.index_word))

Mecab Vocab size : 52279

```

1) tokenizer.sequences_to_texts() 함수를 사용하여 Decoding

2) tokenizer.index_word 를 사용하여 Decoding

두 가지 방법으로 mecab_tensor[100] 을 원문으로 되돌려 보세요! (여기서 띄어쓰기는 고려하지 않습니다!)

```

1 # Case 1 : mecab_tokenizer.sequences_to_texts()
2
3 # Case 1
4 texts = mecab_tokenizer.sequences_to_texts([mecab_tensor[100]]) # 코드 작성
5
6 print(texts[0])

```

중국 관영 언론 은 2000 가지 음식 명 을 제시 한 170 페이지 분량 의 책 을 베이징 호텔 에 제공 했 디

```

1 # Case 2 : mecab_tokenizer.index_word[]
2 sentence = ""
3
4 for w in mecab_tensor[100]:

```

```
5     if w == 0: continue
6     sentence += mecab_tokenizer.index_word[w] + " "# 코드 작성
7
8 print(sentence)
```

중국 관영 언론 은 2000 가지 음식 명 을 제시 한 170 페이지 분량 의 책 을 베이징 호텔 에 제공 했 다

▼ 사전에 없는 단어의 문제

코로나바이러스는 2019년 12월 중국 우한에서 처음 발생한 뒤
전 세계로 확산된, 새로운 유형의 호흡기 감염 질환입니다.

→

<unk>는 2019년 12월 중국 <unk>에서 처음 발생한 뒤
전 세계로 확산된, 새로운 유형의 호흡기 감염 질환입니다.

만약 위 문장을 영문으로 번역해야 한다면 어떨까요? 핵심인 단어 코로나바이러스 와 우한 을 모른다면 제대로 해
낼 수 있을 리가 없습니다. 이를 **OOV(Out-Of-Vocabulary)** 문제라고 합니다. 이처럼 **새로 등장한(본 적 없는)**
단어에 대해 약한 모습을 보일 수밖에 없는 기법들이기에, 이를 해결하고자 하는 시도들이 있었습니다. 그리고 그
것이 우리가 다음 스텝에서 배울, **Wordpiece Model**이죠!

