

Wordpiece Model(WPM) 은 우리가 접한 적이 있는 아이디어를 기반으로 만들어졌습니다. 두 단어 `preview` 와 `predict` 를 보면 접두어인 `pre` 가 공통되고 있죠? `pre` 가 들어간 단어는 주로 "**미리**", "**이전의**" 와 연계되는 의미를 가지고 있습니다. 컴퓨터도 두 단어를 따로 볼 게 아니라 `pre+view` 와 `pre+dict` 로 본다면 학습을 더 잘 할 수 있지 않을까요?

이처럼 한 단어를 여러 개의 Subword의 집합으로 보는 방법이 WPM입니다. WPM의 원리를 알기 전, 먼저 알아야 할 것이 바로 **Byte Pair Encoding(BPE)** 입니다.

Byte Pair Encoding(BPE)

BPE 알고리즘이 고안된 것은 1994년입니다. 그때는 자연어 처리에 적용하기 위해서가 아니라 데이터 압축을 위해서 생겨났었죠. 데이터에서 **가장 많이 등장하는 바이트 쌍(Byte Pair)** 을 새로운 단어로 치환하여 압축하는 작업을 반복하는 방식으로 동작합니다. 예시는 아래와 같습니다.

```
aaabdaaabac # 가장 많이 등장한 바이트 쌍 "aa"를 "z"로 치환합니다.
→
ZabdZabac   # "aa" 총 두 개가 치환되어 4바이트를 2바이트로 압축하였습니다.
Z=aa        # 그다음 많이 등장한 바이트 쌍 "ab"를 "Y"로 치환합니다.
→
ZYdZYac     # "ab" 총 두 개가 치환되어 4바이트를 2바이트로 압축하였습니다.
Z=aa        # 여기서 작업을 멈추어도 되지만, 치환된 바이트에 대해서도 진행한다면
Y=ab        # 가장 많이 등장한 바이트 쌍 "ZY"를 "X"로 치환합니다.
→
XdXac
Z=aa
Y=ab
X=ZY        # 압축이 완료되었습니다!
```

아주 직관적인 알고리즘이죠? 이를 토큰화에 적용하자고 제안한 것은 2015년이었습니다. 모든 단어를 문자(바이트)들의 집합으로 취급하여 자주 등장하는 문자 쌍을 합치면, 접두어나 접미어의 의미를 캐치할 수 있고, 처음 등장하는 단어는 문자(알파벳)들의 조합으로 나타내어 **OOV 문제를 완전히 해결**할 수 있다는 것이죠!

비교적 최근의 기술을 소개해드리는 만큼 논문을 함께 첨부합니다.

- [Neural Machine Translation of Rare Words with Subword Units](#)

위 논문은 Python 소스 코드를 함께 제공해 주어 간편하게 실습을 해 볼 수 있습니다. 논문에서 제공해 주는 예제로 동작 방식을 자세히 들여다보죠!

▼ n-grams

- 유니그램 (1)
- 바이그램 (2)
- 트라이그램 (3)

나는 학교에 가서 공부를 합니다.

```

1 def n_grams(text, n):
2     return [text[i:i+n] for i in range(len(text)-n+1)]
3
4 cleaned = ['mary', ',', 'n't', 'slap', 'green', 'witch', '.']
5 print(n_grams(cleaned, 1))

[['mary'], [','], ['n't'], ['slap'], ['green'], ['witch'], ['.']]

```

▼ BPE Algorithm

```

1 import re, collections

1 num_merges = 10 # BPE를 몇 회 수행할 것지 정함.

```

```

# BPE corpus
corpus = """
low lower newest widest
low lower newest widest
low     newest widest
low     newest
low     newest
        newest
"""

```

```

1 dictionary = {'l o w </w>' : 5,
2              'l o w e r </w>' : 2,
3              'n e w e s t </w>' : 6,
4              'w i d e s t </w>' : 3
5 }

```

```

1 def get_stats(dictionary):
2
3     """
4     단어 사전을 불러와
5     단어는 공백 단위로 쪼개어 문자 list를 만들고
6     빈도수와 쌍을 이루게 합니다. (symbols)
7     """
8
9     # 유니그램의 pair들의 빈도수를 카운트
10    pairs = collections.defaultdict(int)
11    for word, freq in dictionary.items():
12        symbols = word.split()
13        for i in range(len(symbols)-1):

```

```

13         for i in range(len(symbols)-1):
14             pairs[symbols[i], symbols[i+1]] += freq
15     print('현재 pair들의 빈도수 :', dict(pairs))
16     return pairs

```

```

1 def merge_dictionary(pair, v_in):
2     v_out = {}
3     bigram = re.escape(' '.join(pair))
4     p = re.compile(r'(?!\S)' + bigram + r'(?!\S)')
5     for word in v_in:
6         w_out = p.sub(' '.join(pair), word)
7         v_out[w_out] = v_in[word]
8     return v_out

```

```

1 bpe_codes = {}
2 bpe_codes_reverse = {}
3 for i in range(num_merges):
4     print(">> Step {0}".format(i+1))
5     pairs = get_stats(dictionary)
6     best = max(pairs, key=pairs.get)
7     dictionary = merge_dictionary(best, dictionary)
8
9     bpe_codes[best] = i
10    bpe_codes_reverse[best[0] + best[1]] = best
11
12    print("new merge: {}".format(best))
13    print("dictionary: {}".format(dictionary))

```

>> Step 1

현재 pair들의 빈도수 : {('l', 'o'): 7, ('o', 'w'): 7, ('w', '</w>'): 5, ('w', 'e')
new merge: ('e', 's')
dictionary: {'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i

>> Step 2

현재 pair들의 빈도수 : {('l', 'o'): 7, ('o', 'w'): 7, ('w', '</w>'): 5, ('w', 'e')
new merge: ('es', 't')
dictionary: {'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d

>> Step 3

현재 pair들의 빈도수 : {('l', 'o'): 7, ('o', 'w'): 7, ('w', '</w>'): 5, ('w', 'e')
new merge: ('est', '</w>')
dictionary: {'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d

>> Step 4

현재 pair들의 빈도수 : {('l', 'o'): 7, ('o', 'w'): 7, ('w', '</w>'): 5, ('w', 'e')
new merge: ('l', 'o')
dictionary: {'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s

>> Step 5

현재 pair들의 빈도수 : {('lo', 'w'): 7, ('w', '</w>'): 5, ('w', 'e'): 2, ('e', 'r')
new merge: ('lo', 'w')
dictionary: {'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s t <

>> Step 6

현재 pair들의 빈도수 : {('low', '</w>'): 5, ('low', 'e'): 2, ('e', 'r'): 2, ('r',
new merge: ('n', 'e')
dictionary: {'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s t </

>> Step 7

현재 pair들의 빈도수 : {('low', '</w>'): 5, ('low', 'e'): 2, ('e', 'r'): 2, ('r',
new merge: ('ne', 'w')
dictionary: {'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s t </w

```
>> Step 8
현재 pair들의 빈도수 : {('low', '</w>'): 5, ('low', 'e'): 2, ('e', 'r'): 2, ('r',
new merge: ('new', 'est</w>')
dictionary: {'low </w>': 5, 'low e r </w>': 2, 'newest</w>': 6, 'w i d est</w>'}
>> Step 9
현재 pair들의 빈도수 : {('low', '</w>'): 5, ('low', 'e'): 2, ('e', 'r'): 2, ('r',
new merge: ('low', '</w>')
dictionary: {'low</w>': 5, 'low e r </w>': 2, 'newest</w>': 6, 'w i d est</w>'}
>> Step 10
현재 pair들의 빈도수 : {('low', 'e'): 2, ('e', 'r'): 2, ('r', '</w>'): 2, ('w', 'i')
new merge: ('w', 'i')
dictionary: {'low</w>': 5, 'low e r </w>': 2, 'newest</w>': 6, 'wi d est</w>':
```

```
1 print(bpe_codes)
```

```
{('e', 's'): 0, ('es', 't'): 1, ('est', '</w>'): 2, ('l', 'o'): 3, ('lo', 'w')}
```

▼ OOV 대처하기

만일 lowest 라는 처음 보는 단어가 등장하더라도, 위 알고리즘을 따르면 어느 정도 의미가 파악된 low와 est의 결합으로 표현할 수 있습니다. 또 BPE의 놀라운 점은 아무리 큰 데이터도 원하는 크기로 OOV 문제없이 사전을 정의할 수 있다는 것입니다. 극단적으로 생각했을 때 알파벳 26개와 특수문자, 문장부호를 아무리 추가해도 100개 이내로 사전을 정의할 수 있죠.(물론 그러면 안 됩니다!!)

Embedding 레이어는 단어의 개수 x Embedding 차원 수의 Weight를 생성하기 때문에 단어의 개수가 줄어드는 것은 곧 메모리의 절약으로 이어집니다. 많은 데이터가 곧 정확도로 이어지기 때문에 이런 기여는 굉장히 의미가 있습니다!

하지만 아직도! 완벽하다고는 할 수 없습니다. 만약 수많은 데이터를 사용해 만든 BPE 사전으로 모델을 학습시키고 문장을 생성하게 했다고 합시다. 그게 [i, am, a, b, o, y, a, n, d, you, are, a, gir, l] 이라면, 어떤 기준으로 이들을 결합해서 문장을 복원하죠? 몽땅 한꺼번에 합쳤다면 끔찍한 일이 벌어질 것만 같습니다...

```
1 def get_pairs(word):
2     pairs = set()
3     prev_char = word[0]
4     for char in word[1:]:
5         pairs.add((prev_char, char))
6         prev_char = char
7     return pairs
```

```
1 orig = 'hi'
2 word = tuple(orig) + ('</w>',)
3 print(word)

('h', 'i', '</w>')
```

```
1 def encode(orig):
2     word = tuple(orig) + ('</w>',)
```

```

2 word = tuple(orig) + ('</w>',)
3 print("__word split into characters:__ <tt>{}<tt>".format(word))
4
5 pairs = get_pairs(word)
6
7 if not pairs:
8     return orig
9
10 iteration = 0
11 while True:
12     iteration += 1
13     print("__Iteration {}:__".format(iteration))
14
15     print("Bigram in the word: {}".format(pairs))
16     bigram = min(pairs, key = lambda pair: bpe_codes.get(pair, float('inf')))
17     print("candidate for merging: {}".format(bigram))
18     if bigram not in bpe_codes:
19         print("__Candidate not in BPE merges, algorithm stops.__")
20         break
21     first, second = bigram
22     new_word = []
23     i = 0
24     while i < len(word):
25         try:
26             j = word.index(first, i)
27             new_word.extend(word[i:j])
28             i = j
29         except:
30             new_word.extend(word[i:])
31             break
32
33     if word[i] == first and i < len(word)-1 and word[i+1] == second:
34         new_word.append(first+second)
35         i += 2
36     else:
37         new_word.append(word[i])
38         i += 1
39     new_word = tuple(new_word)
40     word = new_word
41     print("word after merging : {}".format(word))
42     if len(word) == 1:
43         break
44     else:
45         pairs = get_pairs(word)
46
47 # 특별토큰인 </w>는 출력하지 않는다.
48 if word[-1] == '</w>':
49     word = word[:-1]
50 elif word[-1].endswith('</w>'):
51     word = word[:-1] + (word[-1].replace('</w>', ''),)
52 return word

```

단어 'loki'가 들어오면 BPE 알고리즘 해당 단어를 어떻게 분리할까요?

```
1 encode("loki")

__word split into characters:__ <tt>('l', 'o', 'k', 'i', '</w>')<tt>
__Iteration 1:__
Bigram in the word: (('o', 'k'), ('l', 'o'), ('i', '</w>'), ('k', 'i'))
candidate for merging: ('l', 'o')
word after merging : ('lo', 'k', 'i', '</w>')
__Iteration 2:__
Bigram in the word: (('i', '</w>'), ('lo', 'k'), ('k', 'i'))
candidate for merging: ('i', '</w>')
__Candidate not in BPE merges, algorithm stops.__
('lo', 'k', 'i')
```

현재 서브워드 단어집합에는 'lo'가 존재하므로, 'lo'는 유지하고 'k'와 'i'는 분리시킵니다. 단어 'lowest'에 대해서도 수행해봅시다.

```
1 encode("lowest")

__word split into characters:__ <tt>('l', 'o', 'w', 'e', 's', 't', '</w>')<tt>
__Iteration 1:__
Bigram in the word: (('t', '</w>'), ('w', 'e'), ('l', 'o'), ('e', 's'), ('s',
candidate for merging: ('e', 's')
word after merging : ('l', 'o', 'w', 'es', 't', '</w>')
__Iteration 2:__
Bigram in the word: (('t', '</w>'), ('l', 'o'), ('w', 'es'), ('es', 't'), ('o'
candidate for merging: ('es', 't')
word after merging : ('l', 'o', 'w', 'est', '</w>')
__Iteration 3:__
Bigram in the word: (('l', 'o'), ('est', '</w>'), ('w', 'est'), ('o', 'w'))
candidate for merging: ('est', '</w>')
word after merging : ('l', 'o', 'w', 'est</w>')
__Iteration 4:__
Bigram in the word: (('l', 'o'), ('o', 'w'), ('w', 'est</w>'))
candidate for merging: ('l', 'o')
word after merging : ('lo', 'w', 'est</w>')
__Iteration 5:__
Bigram in the word: (('w', 'est</w>'), ('lo', 'w'))
candidate for merging: ('lo', 'w')
word after merging : ('low', 'est</w>')
__Iteration 6:__
Bigram in the word: (('low', 'est</w>'))
candidate for merging: ('low', 'est</w>')
__Candidate not in BPE merges, algorithm stops.__
('low', 'est')
```

현재 서브워드 단어집합에 'low'와 'est'가 존재하므로, 'low'와 'est'를 분리시킵니다. 단어 'lowing'에 대해서도 수행해봅시다.

```
1 encode("lowing")

__word split into characters:__ <tt>('l', 'o', 'w', 'i', 'n', 'g', '</w>')<tt>
__Iteration 1:__
Bigram in the word: (('i', 'n'), ('l', 'o'), ('w', 'i'), ('g', '</w>'), ('o',
candidate for merging: ('l', 'o')
```

```
word after merging : ('lo', 'w', 'i', 'n', 'g', '</w>')
__Iteration 2:__
Bigram in the word: {('i', 'n'), ('w', 'i'), ('g', '</w>'), ('lo', 'w'), ('n',
candidate for merging: ('lo', 'w')}
word after merging : ('low', 'i', 'n', 'g', '</w>')
__Iteration 3:__
Bigram in the word: {('n', 'g'), ('g', '</w>'), ('i', 'n'), ('low', 'i')}
candidate for merging: ('n', 'g')
__Candidate not in BPE merges, algorithm stops.__
('low', 'i', 'n', 'g')
```

현재 서브워드 단어집합에 'low'가 존재하지만, 'i', 'n', 'g'의 바이그램 조합으로 이루어진 서브워드는 존재하지 않으므로 'i', 'n', 'g'로 전부 분리합니다. 훈련된 데이터 중에서 어떤 서브워드도 존재하지 않는 'highing'은 어떨까요?

```
1 encode("highing")
```

```
__word split into characters:__ <tt>('h', 'i', 'g', 'h', 'i', 'n', 'g', '</w>')
__Iteration 1:__
Bigram in the word: {('i', 'n'), ('h', 'i'), ('g', '</w>'), ('i', 'g'), ('n',
candidate for merging: ('i', 'n')}
__Candidate not in BPE merges, algorithm stops.__
('h', 'i', 'g', 'h', 'i', 'n', 'g')
```

▼ Wordpiece Model(WPM)

이에 구글에서 BPE를 변형해 제안한 알고리즘이 바로 WPM입니다. WPM은 BPE에 대해 두 가지 차별성을 가집니다.

1. 공백 복원을 위해 단어의 시작 부분에 언더바 _ 를 추가합니다.
2. 빈도수 기반이 아닌 가능도(Likelihood)를 증가시키는 방향으로 문자 쌍을 합칩니다.

첫 번째 문항은 아주 쉬운 내용으로, 앞서 사용한 예문을 빌리면 [_i, _am, _a, _b, o, y, _a, n, d, _you, _are, _a, _gir, l] 로 토큰화를 한다는 것입니다. 이렇게 하면 문장을 복원하는 과정이 **1) 모든 토큰을 합친 후, 2) 언더바 _를 공백으로 치환**으로 마무리되어 간편하죠.

두 번째 문항은 다소 난해하게 다가올 수도 있습니다. 직관적인 이해를 얻고 넘어가는 것을 목표로 하죠. 본 내용은 아래 논문 3절과 4절에 자세하게 나와 있습니다.

- [JAPANESE AND KOREAN VOICE SEARCH](#)

(여기서 잠깐) 구글에서 이 기법을 한국어, 일본어 텍스트 처리를 위해 고려했다는 사실이 흥미롭지 않나요? 이 점은 2가지를 시사합니다.

- 조사, 어미 등의 활용이 많고 복잡한 한국어 같은 모델의 토큰나이저로 WPM이 좋은 대안이 될 수 있다.
- WPM은 어떤 언어든 무관하게 적용 가능한 language-neutral하고 general한 기법이다. 한국어 형태소 분석기처럼 한국어에만 적용 가능한 기법보다 훨씬 활용도가 크다.

WPM을 수행하기 이전의 문장: Jet makers feud over seat width with big orders at stake

WPM을 수행한 결과(wordpieces): _J et _makers _fe ud _over _seat _width _with _big _orders _at _stake

Jet는 J와 et로 나누어졌으며, feud는 fe와 ud로 나누어진 것을 볼 수 있습니다. WPM은 모든 단어의 맨 앞에 _를 붙이고, 단어는 서브 워드(subword)로 통계에 기반하여 띄어쓰기로 분리합니다. 여기서 언더바 _는 문장 복원을 위한 장치입니다.

예컨대, WPM의 결과로 나온 문장을 보면, Jet → _J et와 같이 기존에 없던 띄어쓰기가 추가되어 서브 워드(subwords)들을 구분하는 구분자 역할을 하고 있습니다. 그렇다면 기존에 있던 띄어쓰기와 구분자 역할의 띄어쓰기는 어떻게 구별할까요? 이 역할을 수행하는 것이 단어들 앞에 붙은 언더바 _입니다. WPM이 수행된 결과로부터 다시 수행 전의 결과로 돌리는 방법은 현재 있는 모든 띄어쓰기를 전부 제거하고, 언더바를 띄어쓰기로 바꾸면 됩니다.

이 알고리즘은 최신 딥 러닝 모델 BERT를 훈련하기 위해서 사용되기도 하였습니다.

토큰화의 끝판왕으로 보이는 이 WPM은 아쉽게도 공개되어 있지는 않습니다. 대신에 구글의 *SentencePiece* 라이브러리를 통해 고성능의 BPE를 사용할 수 있습니다! *SentencePiece*에는 전처리 과정도 포함되어 있어서, 데이터를 따로 정제할 필요가 없어 간편하기까지 합니다. 따라서 아래 깃허브 페이지에 방문해 사용법을 간단히 눈으로 봐두도록 하죠! 아마 다음 프로젝트 노트에서는 *SentencePiece* 라이브러리를 적극적으로 활용하는 실습을 진행하게 될 것입니다.

- [google/sentencepiece](https://google.github.io/sentencepiece/)

이제 우리는 어떤 언어에도 OOV 발생 우려 없이 안정적으로 활용할 수 있는 멋진 토큰나이징 기술을 확보했습니다. 이제는 컴퓨터가 단어사전을 안심하고 활용할 수 있겠군요!

하지만 아직 우리는 제대로 된 단어의 분산 표현을 얻는 법을 제대로 다루지는 않았습니다. 이쯤에서 이런 고민이 생기게 됩니다. 한국어라면 자동차를 _자동 / 차 로 분리되는데... 속성이 아무리 추상적이래도 보기에 차가 마시는 차인지, 달리는 차인지 도통 알 수가 없죠? 게다가 설령 토큰화가 완벽하다고 해도, 남자가 [-1, 0] 인지 [1, 0] 인지는 컴퓨터 입장에서는 알 도리가 없습니다.

Embedding 레이어는 선언 즉시 랜덤한 실수로 Weight 값을 채우고, 학습을 진행하며 적당히 튜닝해가는 방식으로 속성을 맞춰가지만 이는 뭔가 찝찝합니다. 토큰들이 멋지게 의미를 갖게 하는 방법은 없을까요?

▶ IMDB리뷰 tensorflow의 subwordTextEncoder로 토큰화

[] ↪ 숨겨진 셀 17개

▶ 네이버 영화 리뷰로 서브워드토큰나이저

[] ↪ 숨겨진 셀 13개

▼ IMDB리뷰 sentencePiece로 토큰화

```
1 !pip install sentencepiece
```

```
Collecting sentencepiece
```


[illegible]

```
1 import sentencepiece as spm
2 import pandas as pd
3 import urllib.request
4 import csv

1 urllib.request.urlretrieve("https://raw.githubusercontent.com/LawrenceDuan/IMDb-
    ('IMDb_Reviews.csv', <http.client.HTTPMessage at 0x7fea90b77d50>)

1 train_df = pd.read_csv('IMDb_Reviews.csv')

1 train_df.head()
```

	review	sentiment
0	My family and I normally do not watch local mo...	1
1	Believe it or not, this was at one time the wo...	0
2	After some internet surfing, I found the "Home...	0
3	One of the most unheralded great works of anim...	1
4	It was the Sixties, and anyone with long hair ...	0

```
1 print('리뷰 개수 : ', len(train_df))
```

리뷰 개수 : 50000

```
1 with open('imdb_review.txt', 'w', encoding='utf8') as f:
2     f.write('\n'.join(train_df['review']))
```

sentencepiece의 자세한 실행 옵션 : <https://github.com/google/sentencepiece>

```
spm.SentencePieceTrainer.train(
    f"--input={corpus} --model_prefix={prefix} --vocab_size={vocab_size + 7}" +
    " --model_type=bpe" +
    " --max_sentence_length=999999" + # 문장 최대 길이
    " --pad_id=0 --pad_piece=[PAD]" + # pad (0)
    " --unk_id=1 --unk_piece=[UNK]" + # unknown (1)
    " --bos_id=2 --bos_piece=[BOS]" + # begin of sequence (2)
    " --eos_id=3 --eos_piece=[EOS]" + # end of sequence (3)
    " --user_defined_symbols=[SEP],[CLS],[MASK]") # 사용자 정의 토큰
```

```
spm.SentencePieceTrainer.Train('--input=imdb_review.txt --model_prefix=imdb --vocab_size=5
```

```
1 corpus = "imdb_review.txt" #입력 corpus
2 prefix = "imdb" # 저장할 단어장 이름
3 vocab_size = 5000
4 spm.SentencePieceTrainer.Train(
5     f"--input={corpus} --model_prefix={prefix} --vocab_size={vocab_size}" +
6     " --model_type=bpe" +
7     " --max_sentence_length=9999") # 문장 최대 길이
```

- input : 학습시킬 파일
- model_prefix : 만들어질 모델 이름
- vocab_size : 단어집합크기
- model_type : 사용할 모델 (unigram(default), bpe, char, word)
- pad_id, pad_piece : pad token id, 값
- unk_id, unk_piece : unknown token id, 값
- bos_id, bos_piece : begin of sentence token id, 값
- eos_id, eos_piece : end of sequence token id, 값
- user_defined_symbols : 사용자 정의 토큰

```
1 vocab_list = pd.read_csv('imdb.vocab', sep='\t', header=None, quoting=csv.QUOTE_
2 vocab_list.sample(10)
```

	0	1
2997	__tort	-2994
1498	__brilliant	-1495
471	__being	-468
2511	__public	-2508
4650	__theatre	-4647
3223	__exception	-3220
4581	John	-4578
3664	aint	-3661
4260	__brutal	-4257
4589	__hadn	-4586

```
1 len(vocab_list)
```

```
1 sp = spm.SentencePieceProcessor()
2 vocab_file = "imdb.model"
3 sp.load(vocab_file)
```

```
True
```

```
1 lines = [
2     "I didn't at all think of it this way.",
3     "I have waited a long time for someone to film"
4 ]
5 for line in lines:
6     print(line)
7     print(sp.encode_as_pieces(line)) # 문장을 입력하면 서브워드 시퀀스로 변환
8     print(sp.encode_as_ids(line)) # 문장을 입력하면 정수 시퀀스로 변환
9     print()

I didn't at all think of it this way.
['_I', '_didn', "'", 't', '_at', '_all', '_think', '_of', '_it', '_this', '_wa
[41, 623, 4950, 4926, 138, 169, 378, 30, 58, 73, 413, 4945]

I have waited a long time for someone to film
['_I', '_have', '_wa', '_ited', '_a', '_long', '_time', '_for', '_someone', '_t
[41, 141, 1364, 1120, 4, 666, 285, 92, 1078, 33, 91]
```

```
1 sp.GetPieceSize() #단어집합의 크기
```

```
5000
```

```
1 sp.IdToPiece(120) #정수로부터 매핑되는 서브워드 변환
```

```
'_not'
```

```
1 # 위에서 출력된 결과를 sp.PieceToId 안으로 입력하기
2 sp.PieceToId('_not') # 서브워드로부터 매핑되는 정수로 변환
```

```
120
```

```
1 # 위에서 출력된 정수 시퀀스를 sp.DecodeIds의 정수 시퀀스에 입력하기
2 sp.DecodeIds([41, 141, 1364, 1120, 4, 666, 285, 92, 1078, 33, 91])
```

```
'I have waited a long time for someone to film'
```

```
1 # 서브워드 시퀀스로부터 문장으로 변환
```

```
2 sp.DecodePieces(['_I', '_have', '_wa', '_ited', '_a', '_long', '_time', '_for', '_
'I have waited a long time for someone to film'
```

```
1 print(sp.encode('I have waited a long time for someone to film', out_type=str))
2 print(sp.encode('I have waited a long time for someone to film', out_type=int))

['_I', '_have', '_wa', '_ited', '_a', '_long', '_time', '_for', '_someone', '_t
[41, 141, 1364, 1120, 4, 666, 285, 92, 1078, 33, 91]
```

▼ 네이버 영화리뷰로 센텐스피스 적용하기

```
1 import pandas as pd
2 import sentencepiece as spm
3 import urllib.request
4 import csv
```

```
1 urllib.request.urlretrieve("https://raw.githubusercontent.com/e9t/nsmc/master/ratings.txt", <http.client.HTTPMessage at 0x7fea92634410>)
```

```
1 naver_df = pd.read_table('ratings.txt')
2 naver_df.head()
```

	id	document	label
0	8112052	어릴때보고 지금다시봐도 재밌어요ㅋㅋ	1
1	8132799	디자인을 배우는 학생으로, 외국디자이너와 그들이 일군 전통을 통해 발전해가는 문화 산...	1
2	4655635	폴리스스토리 시리즈는 1부터 뉴까지 버릴게 하나도 없음.. 최고.	1
3	9251303	와.. 연기가 진짜 개쩔구나.. 지루할거라고 생각했는데 몰입해서 봤다.. 그래 이런...	1
4	10067386	안개 자욱한 밤하늘에 떠 있는 초승달 같은 영화.	1

```
1 # 리뷰 갯수 출력하기
2 print('리뷰 개수 :',len(naver_df)) # 리뷰 개수 출력
```

리뷰 개수 : 200000

```
1 # Null 존재하는지 확인하기
2 print(naver_df.isnull().values.any())
```

True

```
1 # Null값이 존재하는 행 제거
2 naver_df = naver_df.dropna(how = 'any')
3 # Null 존재하는지 확인하기
4 print(naver_df.isnull().values.any())
```

False

```
1 # 리뷰 갯수 출력하기
2 print('리뷰 개수 :',len(naver_df))
```

리뷰 개수 : 199992

```
1 with open('naver_review.txt', 'w', encoding='utf8') as f:
2     f.write('\n'.join(naver_df['document']))
```

```
1 spm.SentencePieceTrainer.Train('--input=naver_review.txt --model_prefix=naver --
```

```
1 # vocab불러오기 pd.read_csv
2 vocab_list = pd.read_csv('naver.vocab', sep='\t', header=None, quoting=csv.QUOTE
3 #vocab_list[:10]
4 vocab_list.head(10)
5 # vocab 10개출력
```

	0	1
0	<unk>	0
1	<s>	0
2	</s>	0
3	..	0
4	영화	-1
5	__영화	-2
6	__이	-3
7	__아	-4
8	...	-5
9	__그	-6

```
1 sp = spm.SentencePieceProcessor()
2 vocab_file = "naver.model"
3 sp.load(vocab_file)
```

True

```
1 lines = [
2     "뭐 이딴 것도 영화냐.",
3     "진짜 최고의 영화입니다 ㅋㅋ",
4 ]
5 for line in lines:
6     print(line)
7     print(sp.encode_as_pieces(line))
8     print(sp.encode_as_ids(line))
9     print()
```

뭐 이딴 것도 영화냐.
 ['_뭐', '_이딴', '_것도', '_영화냐', '._']
 [132, 966, 1296, 2590, 3276]

진짜 최고의 영화입니다 ㅋㅋ
 ['_진짜', '_최고의', '_영화입니다', '._ㅋㅋ']
 [54, 200, 821, 85]

```

1 sp.GetPieceSize()

5000

1 sp.IdToPiece(4)

'영화'

1 sp.PieceToId('영화')

4

1 sp.DecodeIds([132, 966, 1296, 2590, 3276])

'뭐 이딴 것도 영화냐.'

1 sp.DecodePieces(['_뭐', '_이딴', '_것도', '_영화냐', '.']) # 서브워드 시퀀스 입력

'뭐 이딴 것도 영화냐.'

1 print(sp.encode('진짜 최고의 영화입니다 ㅋㅋ', out_type=str))
2 print(sp.encode('진짜 최고의 영화입니다 ㅋㅋ', out_type=int))

['_진짜', '_최고의', '_영화입니다', '_ㅋㅋ']
[54, 200, 821, 85]

```

▼ one-hot encoding실습

```

1 import numpy as np

1 # 입력 문장
2 raw_inputs = [
3     "나는 학생 입니다",
4     "나는 좋은 선생님 입니다",
5     "당신은 매우 좋은 선생님 입니다"
6 ]
7
8 # 정답 학생(1) 기타(0)
9 raw_labels = [1, 0, 0]

1 words = []
2 for s in raw_inputs:
3     words.extend(s.split())
4
5 # 중복 단어 제거
6 words = list(dict.fromkeys(words))
7
8 word_to_id = {"[PAD]":0, "[UNK]":1}

```

```

9 for w in words:
10     word_to_id[w] = len(word_to_id)
11
12 # 각 번호별 단어
13 id_to_word = {i: w for w, i in word_to_id.items()}
14
15 print(id_to_word)

```

```

{0: '[PAD]', 1: '[UNK]', 2: '나는', 3: '학생', 4: '입니다', 5: '좋은', 6: '선생님',

```

```

1 train_inputs = []
2 for s in raw_inputs:
3     row = [word_to_id[w] for w in s.split()]
4     row += [0] * (5- len(row))
5     train_inputs.append(row)
6
7 train_inputs = np.array(train_inputs)
8
9 print(train_inputs)

```

```

[[2 3 4 0 0]
 [2 5 6 4 0]
 [7 8 5 6 4]]

```

```

1 # one-hot matrix 생성
2 onehot_matrix = np.eye(len(word_to_id))
3 print(onehot_matrix)

```

```

[[1. 0. 0. 0. 0. 0. 0. 0. 0.]
 [0. 1. 0. 0. 0. 0. 0. 0. 0.]
 [0. 0. 1. 0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 1. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 1. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 1. 0. 0. 0.]
 [0. 0. 0. 0. 0. 0. 1. 0. 0.]
 [0. 0. 0. 0. 0. 0. 0. 1. 0.]
 [0. 0. 0. 0. 0. 0. 0. 0. 1.]]

```

```

1 train_onehots = onehot_matrix[train_inputs]
2 print(train_onehots)

```

```

[[[0. 0. 1. 0. 0. 0. 0. 0. 0.]
  [0. 0. 0. 1. 0. 0. 0. 0. 0.]
  [0. 0. 0. 0. 1. 0. 0. 0. 0.]
  [1. 0. 0. 0. 0. 0. 0. 0. 0.]
  [1. 0. 0. 0. 0. 0. 0. 0. 0.]]

[[0. 0. 1. 0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 1. 0. 0. 0.]
 [0. 0. 0. 0. 0. 0. 1. 0. 0.]
 [0. 0. 0. 0. 1. 0. 0. 0. 0.]
 [1. 0. 0. 0. 0. 0. 0. 0. 0.]]

[[0. 0. 0. 0. 0. 0. 0. 1. 0.]
 [0. 0. 0. 0. 0. 0. 0. 0. 1.]]

```

```
[0. 0. 0. 0. 0. 1. 0. 0. 0.]
[0. 0. 0. 0. 0. 0. 1. 0. 0.]
[0. 0. 0. 0. 1. 0. 0. 0. 0.]
```

```
1 print(np.argmax(train_onehots, axis=-1))
```

```
[[2 3 4 0 0]
 [2 5 6 4 0]
 [7 8 5 6 4]]
```

```
1 x = np.argmax(train_onehots, axis=-1)
```

```
1 import tensorflow as tf
2 import tensorflow.keras.layers as L
```

```
1 x_len = train_onehots.shape
```

```
1 inp = tf.convert_to_tensor(x, dtype=tf.int32)
2 inp_len = tf.convert_to_tensor(x_len, dtype=tf.int32)
```

```
1 inp, inp_len
```

```
(<tf.Tensor: shape=(3, 5), dtype=int32, numpy=
array([[2, 3, 4, 0, 0],
       [2, 5, 6, 4, 0],
       [7, 8, 5, 6, 4]], dtype=int32)>,
 <tf.Tensor: shape=(3,), dtype=int32, numpy=array([3, 5, 9], dtype=int32)>)
```

```
1 vocab = 1000
2 dim = 3
3 embed = L.Embedding(vocab, dim)
```

```
1 embed(inp)
```

```
<tf.Tensor: shape=(3, 5, 3), dtype=float32, numpy=
array([[[ 0.0399647 ,  0.0347276 ,  0.03682219],
        [-0.00117704, -0.02982651, -0.02315493],
        [-0.01293083,  0.04588192, -0.0193073 ],
        [-0.0154045 ,  0.00741573,  0.02601862],
        [-0.0154045 ,  0.00741573,  0.02601862]],
       [[ 0.0399647 ,  0.0347276 ,  0.03682219],
        [-0.01873317,  0.03784363, -0.00091445],
        [ 0.03667052, -0.00820971,  0.02016015],
        [-0.01293083,  0.04588192, -0.0193073 ],
        [-0.0154045 ,  0.00741573,  0.02601862]],
       [[ 0.03118033,  0.04501085,  0.02484627],
        [-0.006697 ,  0.03062637,  0.03237622],
        [-0.01873317,  0.03784363, -0.00091445],
        [ 0.03667052, -0.00820971,  0.02016015],
        [-0.01293083,  0.04588192, -0.0193073 ]]], dtype=float32)>
```



```

8 token = txt.morphs(token)
9
10 word2index = {}
11 bow = []
12
13 for voca in token:
14     if voca not in word2index.keys():
15         word2index[voca] = len(word2index)
16         # token을 읽으면서, word2index에 없는 단어는 새로 추가하고, 이미 있는 단어는 넘긴다.
17         bow.insert(len(word2index)-1, 1)
18         # bow전체에 전부 기본값 1을 넣어준다. 단어갯수는 최소 1개이상이기 때문에
19     else:
20         index = word2index.get(voca)
21         # 재 등장하는 단어의 인덱스를 받아오기
22         bow[index] = bow[index]+1
23         # 재 등장하는 단어는 해당하는 인덱스의 위치에 1을 더해줌 (단어갯수를 세는 것)
24
25 print(word2index)

{'소비자': 0, '는': 1, '주로': 2, '소비': 3, '하는': 4, '상품': 5, '을': 6, '기준': 7,

1 bow

[1, 1, 1, 1, 1, 1, 2, 1, 1, 1, 1]

```

▼ tensorflow의 keras Tokenizer를 활용한 BOW

```

1 from tensorflow.keras.preprocessing.text import Tokenizer

1 sentence = ["John likes to watch movies. \
2             Mary likes movies too! \
3             Mary also likes to watch football games."]

1 def print_bow(sentence):
2     tokenizer = Tokenizer()
3     tokenizer.fit_on_texts(sentence) #단어장 생성
4     bow = dict(tokenizer.word_counts) # 각 단어와 각 단어의 빈도를 bow를 저장
5
6     print("Bag of words : ", bow)
7     print("단어장(vocabulary)의 크기 :", len(tokenizer.word_counts)) #중복을 제거한 단어

1 print_bow(sentence)

Bag of words : {'john': 1, 'likes': 3, 'to': 2, 'watch': 2, 'movies': 2, 'mar
단어장(vocabulary)의 크기 : 10

```

▼ scikit-learn의 CountVectorizer를 활용한 BoW

```

1 from sklearn.feature_extraction.text import CountVectorizer

```

```

2 sentence = ["John likes to watch movies. Mary likes movies too! Mary also likes
3
4 vector = CountVectorizer()
5
6 print("Bagg of Words : ", vector.fit_transform(sentence).toarray()) #코퍼스로부터 각
7 print("각 단어의 인덱스 :", vector.vocabulary_) #각 단어의 인덱스가 어떻게 부여되는지 보여줌

Bagg of Words :  [[1 1 1 1 3 2 2 2 1 2]]
각 단어의 인덱스 : {'john': 3, 'likes': 4, 'to': 7, 'watch': 9, 'movies': 6, 'mary

```

▼ 불용어를 제거한 BOW 만들기

▼ 사용자가 직접 정의한 불용어 사용

```

1 from sklearn.feature_extraction.text import CountVectorizer
2
3 text = ["Family is not an important thing. It's everything."]
4
5 vect = CountVectorizer(stop_words=["the", "a", "an", "is", "not"])
6 print(vect.fit_transform(text).toarray())
7 print(vect.vocabulary_)

[[1 1 1 1 1]]
{'family': 1, 'important': 2, 'thing': 4, 'it': 3, 'everything': 0}

```

▼ CountVectorizer에서 제공하는 자체 불용어 사용

```

1 from sklearn.feature_extraction.text import CountVectorizer
2
3 text = ["Family is not an important thing. It's everything."]
4
5 vect = CountVectorizer(stop_words="english")
6 print(vect.fit_transform(text).toarray())
7 print(vect.vocabulary_)

[[1 1 1]]
{'family': 0, 'important': 1, 'thing': 2}

```

▼ NLTK에서 지원하는 불용어 사용

```

1 !pip install nltk
2 import nltk
3 nltk.download('stopwords')

Requirement already satisfied: nltk in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (
[nltk_data] Downloading package stopwords to /root/nltk_data...

```

```
[nltk_data] Unzipping corpora/stopwords.zip.
True
```

```
1 from sklearn.feature_extraction.text import CountVectorizer
2 from nltk.corpus import stopwords
3
4 text = ["Family is not an important thing. It's everything."]
5
6 sw = stopwords.words("english")
7 vect = CountVectorizer(stop_words= sw)
8 print(vect.fit_transform(text).toarray())
9 print(vect.vocabulary_)

[[1 1 1 1]]
{'family': 1, 'important': 2, 'thing': 3, 'everything': 0}
```

▼ DTM (Document-Term Matrix)

다수의 문서에서 등장하는 각 단어들의 빈도를 행렬로 표현한 것
다수의 문서에 대해서 Bow를 하나의 행렬로 표현하고 부르는 용어

문서 1 : I like dog

문서 2 : I like cat

문서 3 : I like cat I like cat

```
1 import pandas as pd
2 content = [[0, 1, 1, 1], [1, 0, 1, 1], [2, 0, 2, 2]]
3 df = pd.DataFrame(content)
4 df.index = ['(문서1) I like dog', '(문서2) I like cat', '(문서3) I like cat I like cat']
5 df.columns = ['cat', 'dog', 'I', 'like']
6 df
```

	cat	dog	I	like
(문서1) I like dog	0	1	1	1
(문서2) I like cat	1	0	1	1
(문서3) I like cat I like cat	2	0	2	2

```
1 import numpy as np
2 from numpy import dot
3 from numpy.linalg import norm
4
5 doc1 = np.array([0, 1, 1, 1])
6 doc2 = np.array([1, 0, 1, 1])
7 doc3 = np.array([2, 0, 2, 2])
8
9 def cos_sim(A, B): # 코사인 유사도는 0~1사이의 값을 가지고, 1에 가까울수록 유사도 높다고 판단
```

```
10     return dot(A,B)/(norm(A) * norm(B))
```

```
1 print(cos_sim(doc1, doc2))
2 print(cos_sim(doc1, doc3))
3 print(cos_sim(doc2, doc3))
```

```
0.6666666666666667
0.6666666666666667
1.0000000000000002
```

▼ scikit-learn CountVectorizer를 활용한 DTM구현

```
1 from sklearn.feature_extraction.text import CountVectorizer
2
3 corpus = [
4     'John likes to watch movies',
5     'Mary likes movies too',
6     'Mary also likes to watch football games',
7 ]
8
9 vector = CountVectorizer()
10 print(vector.fit_transform(corpus).toarray())
11 print(vector.vocabulary_)

[[0 0 0 1 1 0 1 1 0 1]
 [0 0 0 0 1 1 1 0 1 0]
 [1 1 1 0 1 1 0 1 0 1]]
{'john': 3, 'likes': 4, 'to': 7, 'watch': 9, 'movies': 6, 'mary': 5, 'too': 8,
```

▼ TF-IDF (Term-Frequency-Inverse Document Frequency)

모든 문서에서 자주 등장하는 단어는 중요도가 낮다고 판단하고, 특정 문서에서만 자주 등장하는 단어는 중요도가 높다고 판단하는 것

```
1 from math import log
2 import pandas as pd
3
4 docs = [
5     'John likes to watch movies and Mary likes movies too',
6     'James likes to watch TV',
7     'Mary also likes to watch football games',
8 ]

1 vocab = list(set(w for doc in docs for w in doc.split()))
2 vocab.sort()
3
4 print('단어장의 크기 :', len(vocab))
5 print(vocab)
```

단어장의 크기 : 13

['James', 'John', 'Mary', 'TV', 'also', 'and', 'football', 'games', 'likes', 'to', 'too', 'was']

```
1 N = len(docs)
```

```
2 N
```

```
3
```

$$idf(d, t) = \log \frac{n}{1 + df(t)}$$

```
1 def tf(t, d): # 특정 문서 d에서의 특정 단어 t의 등장 횟수
```

```
2     return d.count(t)
```

```
3
```

```
4 def idf(t): # 반비례하는 수
```

```
5     df = 0 #특정 단어 t가 등장한 문서의 수
```

```
6     for doc in docs:
```

```
7         df += t in doc
```

```
8     return log(N/(df+1))+1
```

```
9
```

```
10 def tfidf(t, d):
```

```
11     return tf(t, d) * idf(t)
```

```
1 result = []
```

```
2 for i in range(N):
```

```
3     result.append([])
```

```
4     d = docs[i]
```

```
5     for j in range(len(vocab)):
```

```
6         t = vocab[j]
```

```
7
```

```
8         result[-1].append(tf(t, d))
```

```
9
```

```
10 tf_ = pd.DataFrame(result, columns=vocab)
```

```
11 tf_
```

	James	John	Mary	TV	also	and	football	games	likes	movies	to	too	was
0	0	1	1	0	0	1	0	0	2	2	2	1	
1	1	0	0	1	0	0	0	0	1	0	1	0	
2	0	0	1	0	1	0	1	1	1	0	1	0	

```
1 result = []
```

```
2 for j in range(len(vocab)):
```

```
3     t = vocab[j]
```

```
4     result.append(idf(t))
```

```
5
```

```
6 idf_ = pd.DataFrame(result, index=vocab, columns=['IDF'])
```

```
7 idf_
```



	IDF
James	1.405465
John	1.405465
Mary	1.000000
TV	1.405465
also	1.405465
and	1.405465
football	1.405465
games	1.405465
likes	0.712318
movies	1.405465
to	0.712318
too	1.405465
watch	0.712318

```
1 result = []
2 for i in range(N):
3     result.append([])
4     d = docs[i]
5     for j in range(len(vocab)):
6         t = vocab[j]
7
8         result[-1].append(tfidf(t,d))
9 tfidf_=pd.DataFrame(result, columns=vocab)
10 tfidf_
```

	James	John	Mary	TV	also	and	football	games	likes
0	0.000000	1.405465	1.0	0.000000	0.000000	1.405465	0.000000	0.000000	1.424636
1	1.405465	0.000000	0.0	1.405465	0.000000	0.000000	0.000000	0.000000	0.712318
2	0.000000	0.000000	1.0	0.000000	1.405465	0.000000	1.405465	1.405465	0.712318

```
'0 : John likes to watch movies and Mary likes movies too',
'1 :James likes to watch TV',
'2 : Mary also likes to watch football games',
```

