



+ 코드 + 텍스트

✓ RAM
디스크

BOW 직접 구현

✓
0초

```
[1] 1 docs = ['오늘 동물원에서 코끼리 원숭이를 보고 코끼리 원숭이에게 먹이를 줬어',  
2         '오늘 동물원에서 원숭이에게 사과를 줬어']  
3  
4 # '오늘 동물원에서 코끼리 원숭이를 보고 코끼리 원숭이에게 먹이를 줬어 오늘 동물원에서 원숭이에게 사과를 줬어']
```

✓
0초

```
▶ 1 doc_ls = []  
2 for doc in docs:  
3     doc_ls.append(doc.split(' '))  
4 doc_ls
```

```
➦ [['오늘', '동물원에서', '코끼리', '원숭이를', '보고', '코끼리', '원숭이에게', '먹이를', '줬어'],  
   ['오늘', '동물원에서', '원숭이에게', '사과를', '줬어']]
```

1. 유니크한 토큰 사전 구하기

✓
0초

```
[3] 1 from collections import defaultdict  
2  
3 word2id = defaultdict(lambda : len(word2id))  
4 # word2id  
5  
6 for doc in doc_ls:  
7     # print(doc)  
8     for token in doc:  
9         word2id[token]  
10    # print(token)  
11 word2id
```

0초 [3] defaultdict(<function __main__.<lambda>()>, {
 '오늘': 0,
 '동물원에서': 1,
 '코끼리': 2,
 '원숭이를': 3,
 '보고': 4,
 '원숭이에게': 5,
 '먹이를': 6,
 '줬어': 7,
 '사과를': 8})

2. BOW 구하기

0초 [4] 1 import numpy as np
2 BoW_ls = []
3 for i, doc in enumerate(doc_ls):
4 bow = np.zeros(len(word2id), dtype = int)
5 for token in doc:
6 bow[word2id[token]] += 1
7 BoW_ls.append(bow.tolist())
8 # print(BoW_ls)

0초 [5] 1 BoW_ls

[[1, 1, 2, 1, 1, 1, 1, 1, 0], [1, 1, 0, 0, 0, 1, 0, 1, 1]]



```
1 from IPython.core import display as ICD
2 import pandas as pd
3 sorted_vocab = sorted((value, key) for key, value in word2id.items())
4 print('sorted_vocab',sorted_vocab)
5
6 vocab = []
7 for v in sorted_vocab:
8     vocab.append(v[1])
9 print('vocab',vocab)
10 for i in range(len(docs)) :
11     print("문서{} : {}".format(i, docs[i]))
12     ICD.display(pd.DataFrame([BoW_ls[i]], columns=vocab))
13     print("\n\n")
```



sorted_vocab [(0, '오늘'), (1, '동물원에서'), (2, '코끼리'), (3, '원숭이를'), (4, '보고'), (5, '원숭이에게'), (6, '먹이를'), (7, '줬어'), (8, '사과를')]

vocab ['오늘', '동물원에서', '코끼리', '원숭이를', '보고', '원숭이에게', '먹이를', '줬어', '사과를']

문서0 : 오늘 동물원에서 코끼리 원숭이를 보고 코끼리 원숭이에게 먹이를 줬어

오늘 동물원에서 코끼리 원숭이를 보고 원숭이에게 먹이를 줬어 사과를



0	1	1	2	1	1	1	1	1	0
---	---	---	---	---	---	---	---	---	---

문서1 : 오늘 동물원에서 원숭이에게 사과를 줬어

오늘 동물원에서 코끼리 원숭이를 보고 원숭이에게 먹이를 줬어 사과를

0	1	1	0	0	0	1	0	1	1
---	---	---	---	---	---	---	---	---	---

{x}

▼ sklearn

✓
0초

```
[13] 1 docs = ['오늘 동물원에서 코끼리 원숭이를 보고 코끼리 원숭이에게 먹이를 줬어',  
2          '오늘 동물원에서 코끼리 원숭이를 보고 코끼리 원숭이에게 먹이를 줬어 오늘 동물원에서 원숭이에게 사과를 줬어']
```

✓
0초

```
[15] 1 # 토큰빈도계산 : CountVectorizer  
2 from sklearn.feature_extraction.text import CountVectorizer  
3 #선언  
4 count_vect = CountVectorizer()  
5  
6 BoW = count_vect.fit_transform(docs)  
7 BoW.toarray()
```

```
array([[1, 1, 1, 0, 1, 1, 1, 1, 2],  
       [2, 1, 1, 1, 2, 1, 2, 2, 2]])
```