



+ 코드 + 텍스트

✓ RAM 디스크



## ▼ TDM 직접구현



0초

```
[1] 1 docs = ['동물원 코끼리',
2          '동물원 원숭이 바나나',
3          '엄마 코끼리 아기 코끼리',
4          '원숭이 바나나 코끼리 바나나']
```

1. 토큰화 : 공백으로 토큰화 진행



0초

```
[2] 1 doc_ls=[]
2 for doc in docs:
3     doc_ls.append(doc.split(' '))
4 doc_ls
5
6 # doc_ls = [doc.split() for doc in docs] # 위의 for문을 한줄로 쓰는 경우
```

```
[['동물원', '코끼리'],
 ['동물원', '원숭이', '바나나'],
 ['엄마', '코끼리', '아기', '코끼리'],
 ['원숭이', '바나나', '코끼리', '바나나']]
```

2. 유니크한 토큰 사전 구하기

2-1. 빈 딕셔너리 생성



0초

```
[3] 1 from collections import defaultdict
2
3 word2id = defaultdict(lambda : len(word2id))
4 word2id
```

```
defaultdict(<function __main__.<lambda>()>, {})
```

2-2. 위에서 생성한 빈 딕셔너리에 유니크한 토큰 넣기



0초

```
[4] 1 for doc in doc_ls:
2     for token in doc:
3         word2id[token]
4         print(token)
5         print('\t',word2id)
6 word2id
7
8 # [word2id[token] for doc in doc_ls for token in doc ] # 위의 for문을 한줄로 쓰는 경우
```

0초

[4]

동물원

코끼리

동물원

원숭이

바나나

엄마

코끼리

아기

코끼리

원숭이

바나나

코끼리

바나나

defaultdict(<function <lambda> at 0x7cc25fe4ca60>, {'동물원': 0})

defaultdict(<function <lambda> at 0x7cc25fe4ca60>, {'동물원': 0, '코끼리': 1})

defaultdict(<function <lambda> at 0x7cc25fe4ca60>, {'동물원': 0, '코끼리': 1})

defaultdict(<function <lambda> at 0x7cc25fe4ca60>, {'동물원': 0, '코끼리': 1, '원숭이': 2})

defaultdict(<function <lambda> at 0x7cc25fe4ca60>, {'동물원': 0, '코끼리': 1, '원숭이': 2, '바나나': 3})

defaultdict(<function <lambda> at 0x7cc25fe4ca60>, {'동물원': 0, '코끼리': 1, '원숭이': 2, '바나나': 3, '엄마': 4})

defaultdict(<function <lambda> at 0x7cc25fe4ca60>, {'동물원': 0, '코끼리': 1, '원숭이': 2, '바나나': 3, '엄마': 4})

defaultdict(<function <lambda> at 0x7cc25fe4ca60>, {'동물원': 0, '코끼리': 1, '원숭이': 2, '바나나': 3, '엄마': 4, '아기': 5})

defaultdict(<function <lambda> at 0x7cc25fe4ca60>, {'동물원': 0, '코끼리': 1, '원숭이': 2, '바나나': 3, '엄마': 4, '아기': 5})

defaultdict(<function <lambda> at 0x7cc25fe4ca60>, {'동물원': 0, '코끼리': 1, '원숭이': 2, '바나나': 3, '엄마': 4, '아기': 5})

defaultdict(<function <lambda> at 0x7cc25fe4ca60>, {'동물원': 0, '코끼리': 1, '원숭이': 2, '바나나': 3, '엄마': 4, '아기': 5})

defaultdict(<function <lambda> at 0x7cc25fe4ca60>, {'동물원': 0, '코끼리': 1, '원숭이': 2, '바나나': 3, '엄마': 4, '아기': 5})

defaultdict(<function <lambda> at 0x7cc25fe4ca60>, {'동물원': 0, '코끼리': 1, '원숭이': 2, '바나나': 3, '엄마': 4, '아기': 5})

defaultdict(<function \_\_main\_\_.<lambda>()),  
{'동물원': 0, '코끼리': 1, '원숭이': 2, '바나나': 3, '엄마': 4, '아기': 5})

0초

[5]

1 import numpy as np

2

3 TDM = np.zeros((len(word2id), len(doc\_ls)), dtype=int)

4

5 for i, doc in enumerate(doc\_ls):

6 print(doc)

7 for token in doc:

8 TDM[word2id[token], i] += 1 # 해당 토큰의 위치(column)

9 print(token)

10 print(TDM)

11 print('\t')

['동물원', '코끼리']

동물원

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

코끼리

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

['동물원', '원숭이', '바나나']

동물원

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

$$\begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

엄마

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

$$\begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$
$$\begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$
$$\begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 2 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

```
원숭이
[[1 1 0 0]
 [1 0 2 0]
 [0 1 0 1]
 [0 1 0 0]
 [0 0 1 0]
 [0 0 1 0]]
```



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

```
[7] 1 import pandas as pd
    2
    3 doc_names = ['문서'+ str(i) for i in range(len(doc_ls))]
    4 print('doc_names',doc_names)
    5 sorted_vocab = sorted((value, key) for key, value in word2id.items())
    6 vocab = [v[1] for v in sorted_vocab]
    7 df_TDM = pd.DataFrame(TDM, columns=doc_names)
    8 df_TDM['단어'] = vocab
    9 df_TDM.set_index('단어')
   10
```

doc\_names ['문서0', '문서1', '문서2', '문서3']

	문서0	문서1	문서2	문서3
단어				
동물원	1	1	0	0
코끼리	1	0	2	1
원숭이	0	1	0	1
바나나	0	1	0	2
엄마	0	0	1	0
아기	0	0	1	0

## sklearn

```
[8] 1 docs = ['동물원 코끼리',
    2        '동물원 원숭이 바나나',
    3        '엄마 코끼리 아기 코끼리',
    4        '원숭이 바나나 코끼리 바나나']
```

```
[9] 1 from sklearn.feature_extraction.text import CountVectorizer
    2
    3 count_vect = CountVectorizer() # 참고 sklearn은 DTM으로 만들어지게 설정되어 있음.
    4 DTM = count_vect.fit_transform(docs)
    5 DTM.toarray()
```

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

✓  
0초

[10] 1 DTM.toarray().T

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

↑ ↓ ↺ 💬 ⚙️ 📄 🗑️ ⋮

✓  
0초



```
1 import pandas as pd
2
3 doc_names = ['문서'+ str(i) for i in range(len(doc_ls))]
4 vocab = count_vect.get_feature_names_out()
5 print(vocab)
6 df_TDM = pd.DataFrame(DTM.toarray().T, columns=doc_names)
7 df_TDM['단어'] = vocab
8 df_TDM.set_index('단어')
```

['동물원' '바나나' '아기' '엄마' '원숭이' '코끼리']

문서0 문서1 문서2 문서3



단어



동물원	1	1	0	0
바나나	0	1	0	2
아기	0	0	1	0
엄마	0	0	1	0
원숭이	0	1	0	1
코끼리	1	0	2	1