

+ 코드 + 텍스트

연결 ^

↑ ↓ ↺ 💬 ✎ 📄 🗑️ ⋮

TF-IDF 직접구현

```
[ ] 1 docs = ['오늘 동물원에서 원숭이와 코끼리를 봤어',  
2          '동물원에서 원숭이에게 바나나를 줬어 바나나를']
```

1. 토큰화 : 공백으로 토큰화 진행

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

```
⇒ [['오늘', '동물원에서', '원숭이와', '코끼리를', '봤어'],  
    ['동물원에서', '원숭이에게', '바나나를', '줬어', '바나나를']]
```

2. 유니크한 토큰 사전 구하기

2-1. 빈 딕셔너리 생성

```
[ ] 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

```
[ ] 1 for doc in doc_ls:  
2     for token in doc:  
3         word2id[token]  
4         print(token)  
5         print('\t',word2id)  
6 word2id
```



```

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')

```

➞ ['오늘', '동물원에서', '원숭이와', '코끼리를', '봤어']

오늘

```

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

```

동물원에서

```

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

```

원숭이와

```

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

```

```

1 # TF 계산 특정단어등장빈도/문서내 전체등장단어빈도
2 def computeTF(TDM):
3     doc_len = len(docs) # 문서갯수 2개
4     word_len = len(word2id)
5     print('문서수 : ',doc_len)
6     print('전체 단어수 : ',word_len)
7
8     tf = np.zeros((word_len,doc_len))
9     print(tf)
10
11     for doc_i in range(len(doc_ls)):
12         print(docs[doc_i])
13         for word_i in range(len(word2id)) :
14             tf[word_i,doc_i] = TDM[word_i,doc_i]/TDM[:,doc_i].sum()
15         print(tf)
16     return tf

```

```
[ ] 1 tf = computeTF(TDM)
```

문서수 : 2

전체 단어수 : 8

```

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

```

오늘 동물원에서 원숭이와 코끼리를 봤어

```

[[0.2 0. ]
 [0.2 0. ]
 [0.2 0. ]
 [0.2 0. ]
 [0.2 0. ]
 [0. 0. ]
 [0. 0. ]
 [0. 0. ]]

```

동물원에서 원숭이에게 바나나를 줬어 바나나를

```

[[0.2 0. ]
 [0.2 0.2]
 [0.2 0. ]
 [0.2 0. ]
 [0.2 0. ]

```

[] 1 tf

```
array([[0.2, 0. ],
       [0.2, 0.2],
       [0.2, 0. ],
       [0.2, 0. ],
       [0.2, 0. ],
       [0. , 0.2],
       [0. , 0.4],
       [0. , 0.2]])
```

4. IDF 계산하기

IDF = $-\log(\text{단어가 등장한 문서수} / \text{총 문서수})$

```
1 import math
2 # IDF계산 :  $-\log(\text{단어가 등장한 문서수} / \text{총 문서수})$ 
3 def computeIDF(TDM):
4     doc_len = len(docs) # 문서갯수 2개
5     word_len = len(word2id)
6     print('문서수 : ', doc_len)
7     print('전체 단어수 : ', word_len)
8
9     idf = np.zeros(word_len)
10    print(idf)
11
12    for doc_i in range(len(doc_ls)):
13        for word_i in range(len(word2id)) :
14            idf[word_i] = -math.log10(np.count_nonzero(TDM[word_i,:]) / doc_len)
15    print(idf)
16    return idf
```



```

[0.2 0. ]
[0.2 0. ]
[0.2 0. ]
[0.  0. ]
[0.  0. ]
[0.  0. ]]
동물원에서 원숭이에게 바나나를 줬어 바나나를
[[0.2 0. ]
 [0.2 0.2]
 [0.2 0. ]
 [0.2 0. ]
 [0.2 0. ]
 [0.  0.2]
 [0.  0.4]
 [0.  0.2]]
문서수 : 2
전체 단어수 : 8
[0. 0. 0. 0. 0. 0. 0. 0.]
[ 0.30103 -0.          0.30103 0.30103 0.30103 0.30103 0.30103 0.30103]
오늘 동물원에서 원숭이와 코끼리를 봤어
동물원에서 원숭이에게 바나나를 줬어 바나나를
TF-IDF 최종
[[ 0.060206  0.          ]
 [-0.          -0.          ]
 [ 0.060206  0.          ]
 [ 0.060206  0.          ]
 [ 0.060206  0.          ]
 [ 0.          0.060206]
 [ 0.          0.120412]
 [ 0.          0.060206]]

```

```
[ ] 1 tfidf
```

```

array([[ 0.060206,  0.          ],
       [-0.          , -0.          ],
       [ 0.060206,  0.          ],
       [ 0.060206,  0.          ],
       [ 0.060206,  0.          ],
       [ 0.          , 0.060206],
       [ 0.          , 0.120412],
       [ 0.          , 0.060206]])

```

<> 6. 시각화로 확인

6. 시각화로 확인

```
1 import pandas as pd
2
3 sorted_vocab = sorted((value, key) for key, value in word2id.items())
4 vocab = [v[1] for v in sorted_vocab]
5 print('vocab : ', vocab)
6 tfidf = computeTFIDF(TDM)
7 pd.DataFrame(tfidf.T, columns=vocab)
```

⇒ vocab : ['오늘', '동물원에서', '원숭이와', '코끼리를', '봤어', '원숭이에게', '바나나를', '줬어']

문서수 : 2

전체 단어수 : 8

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

오늘 동물원에서 원숭이와 코끼리를 봤어

```
[[0.2 0. ]
 [0.2 0. ]
 [0.2 0. ]
 [0.2 0. ]
 [0.2 0. ]
 [0.  0. ]
 [0.  0. ]
 [0.  0. ]]
```

동물원에서 원숭이에게 바나나를 줬어 바나나를

```
[[0.2 0. ]
 [0.2 0.2]
 [0.2 0. ]
 [0.2 0. ]
 [0.2 0. ]
 [0.  0.2]
 [0.  0.4]
 [0.  0.2]]
```

문서수 : 2

전체 단어수 : 8


```
[ 0.30103 -0.          0.30103  0.30103  0.30103  0.30103  0.30103  0.30103]
```

오늘 동물원에서 원숭이와 코끼리를 봤어
동물원에서 원숭이에게 바나나를 줬어 바나나를
TF-IDF 최종

```
[[ 0.060206  0.          ]  
 [-0.         -0.         ]  
 [ 0.060206  0.          ]  
 [ 0.060206  0.          ]  
 [ 0.060206  0.          ]  
 [ 0.         0.060206 ]  
 [ 0.         0.120412 ]  
 [ 0.         0.060206 ]]
```

| | 오늘 | 동물원에서 | 원숭이와 | 코끼리를 | 봤어 | 원숭이에게 | 바나나를 | 줬어 |
|---|----------|-------|----------|----------|----------|----------|----------|----------|
| 0 | 0.060206 | -0.0 | 0.060206 | 0.060206 | 0.060206 | 0.000000 | 0.000000 | 0.000000 |
| 1 | 0.000000 | -0.0 | 0.000000 | 0.000000 | 0.000000 | 0.060206 | 0.120412 | 0.060206 |

▼ sklearn

```
[ ] 1 docs = ['오늘 동물원에서 원숭이와 코끼리를 봤어',  
2          '동물원에서 원숭이에게 바나나를 줬어 바나나를']  
3  
4 from sklearn.feature_extraction.text import TfidfVectorizer  
5 tfidf = TfidfVectorizer()  
6 tfidf = tfidf.fit(docs)  
7 tfidf.transform(docs).toarray()  
8 vocab = tfidf.get_feature_names_out()
```

```
▶ 1 import pandas as pd  
2 df = pd.DataFrame(tfidf.transform(docs).toarray(), columns = vocab)  
3 df
```

| | 동물원에서 | 바나나를 | 봤어 | 오늘 | 원숭이에게 | 원숭이와 | 줬어 | 코끼리를 |
|---|----------|----------|----------|----------|----------|----------|----------|----------|
| 0 | 0.335176 | 0.000000 | 0.471078 | 0.471078 | 0.000000 | 0.471078 | 0.000000 | 0.471078 |
| 1 | 0.278943 | 0.784088 | 0.000000 | 0.000000 | 0.392044 | 0.000000 | 0.392044 | 0.000000 |