

▼ 1 잠재의미분석 (LSA)

▼ 1.1 실습 템플릿

```
1 class LSA :
2     def __init__(self, doc_ls, topic_num):
3         pass
4
5     # tdm matrix 생성
6     def TDM(self, doc_ls):
7         pass
8
9     # tdm matrix 특이값 분해(SVD)
10    # U, s, Vt로 분해
11    def SVD(self, tdm):
12        pass
13
14    # 토픽별 주요 키워드 출력
15    def TopicModeling(self) :
16        pass
17
18    # 단어 벡터 행렬 생성 dot(U,s)
19    def TermVectorMatrix(self, u, s):
20        pass
21
22    # 문서 벡터 행렬 생성 dot(s,Vt).T
23    def DocVectorMatrix(self, s, vt):
24        pass
25
26    # 키워드를 입력했을 때 단어 벡터 반환
27    def GetTermVector(self, term):
28        pass
29
30    # 문서를 입력했을 때 문서 벡터 반환
31    def GetDocVector(self, doc):
32        pass
33
34    # 단어-문서 벡터 행렬 생성
35    def TermDocVectorMatrix(self, u, s, vt):
36        pass
37
38    # 단어 벡터 행렬에서 단어 간 코사인 유사도 측정하여 행렬형태로 반환
```

```

39 def TermSimilarityMatrix(self, term_vec_matrix):
40     pass
41
42     # 두개 단어를 입력했을 때 코사인 유사도 반환
43 def GetTermSimilarity(self, term1, term2):
44     pass
45
46     # 문서 벡터 행렬에서 문서 간 코사인 유사도 측정하여 행렬형태로 반환
47 def DocSimilarityMartrix(self, doc_vec_matrix):
48     pass
49
50     # 두개 문서를 입력했을 때 코사인 유사도 반환
51 def GetDocSimilarity(self, doc1, doc2):
52     pass

```

```

1 doc_ls = ['바나나 사과 포도 포도',
2           '사과 포도',
3           '포도 바나나',
4           '짜장면 짬뽕 탕수육',
5           '볶음밥 탕수육',
6           '짜장면 짬뽕',
7           '라면 스시',
8           '스시',
9           '가츠동 스시 소바',
10          '된장찌개 김치찌개 김치',
11          '김치 된장',
12          '비빔밥 김치'
13         ]
14
15 isa = LSA(doc_ls, 3)
16 isa.TopicModeling()
17 isa.GetTermSimilarity('사과', '바나나')
18 isa.GetTermSimilarity('사과', '짜장면')
19 isa.GetDocSimilarity('사과 포도', '포도 바나나')
20 isa.GetDocSimilarity('사과 포도', '라면 스시')

```

▼ 1.2 실습 예제 코드

```

1 import numpy as np
2 from collections import defaultdict
3 from sklearn.metrics.pairwise import cosine_similarity
4 from sklearn.decomposition import randomized_svd
5
6 class LSA :
7     def __init__(self, doc_ls, topic_num):

```

```

8     self.doc_ls = doc_ls
9     self.topic_num = topic_num
10    self.term2idx, self.idx2term = self.toldxDict(' '.join(doc_ls).split())
11    self.doc2idx, self.idx2doc = self.toldxDict(doc_ls)
12
13    self.tdm = self.TDM(doc_ls)
14    self.U, self.s, self.VT = self.SVD(self.tdm)
15
16    #print(self.s)
17    #print(self.U[:, :1].round(2))
18
19    self.term_mat = self.TermVectorMatrix(self.U, self.s, topic_num)
20    self.doc_mat = self.DocVectorMatrix(self.s, self.VT, topic_num)
21    self.term_doc_mat = self.TermDocVectorMatrix(self.U, self.s, self.VT, topic_num)
22
23    #print(self.term2idx.keys())
24    #print(self.term_mat.round(2))
25
26    self.term_sim = self.TermSimilarityMatrix(self.term_mat)
27    self.doc_sim = self.DocSimilarityMartrix(self.doc_mat)
28
29    # 리스트내 값을 index로 변환하는 dict과
30    # index를 리스트내 값으로 변환하는 dict
31    def toldxDict(self, ls) :
32        any2idx = defaultdict(lambda : len(any2idx))
33        idx2any = defaultdict()
34
35        for item in ls:
36            any2idx[item]
37            idx2any[any2idx[item]] = item
38
39        return any2idx, idx2any
40
41    def TDM(self, doc_ls):
42        # 행(토큰크기), 열(문서갯수)로 TDM 생성
43        tdm = np.zeros([len(self.term2idx.keys()), len(doc_ls)])
44
45        for doc_idx, doc in enumerate(doc_ls) :
46            for term in doc.split() :
47                #등장한 단어를 dictionary에서 위치를 탐색하여 빈도수 세기
48                tdm[self.term2idx[term], doc_idx] += 1
49
50        return tdm
51
52    # 특이값 분해
53    def SVD(self, tdm):
54        U, s, VT = randomized_svd(X,
55                                   n_components=15,
56                                   n_iter=5,
57                                   random_state=None)
58

```

```

59     #U, s, VT = np.linalg.svd(tdm)
60     return U, s, VT
61
62     # 토픽별 주요 키워드 출력
63     def TopicModeling(self, count = 3) :
64         topic_num = self.topic_num
65
66         for i in range(topic_num) :
67             score = self.U[:,i:i+1].T
68             sorted_index = np.flip(np.argsort(-score),0)
69
70             a = []
71             for j in sorted_index[0,: count] :
72                 a.append((self.idx2term[j], score[0,j].round(3)))
73
74             print("Topic {} - {}".format(i+1,a ))
75
76     def vectorSimilarity(self, matrix) :
77         similarity = np.zeros([matrix.shape[1], matrix.shape[1]])
78
79         for i in range(matrix.shape[1]) :
80             for j in range(matrix.shape[1]) :
81                 similarity[i,j] = cosine_similarity(matrix[:,i].T, matrix[:,j].T)
82
83         return similarity
84
85     # 키워드를 입력했을 때 단어 벡터 반환
86     def GetTermVector(self, term):
87         vec = self.term_mat[self.term2idx[term]:self.term2idx[term]+1,:]
88         print('{} = {}'.format(term, vec))
89         return vec
90
91     # 문서를 입력했을 때 문서 벡터 반환
92     def GetDocVector(self, doc):
93         vec = self.doc_mat.T[self.doc2idx[doc]:self.doc2idx[doc]+1,:]
94         print('{} = {}'.format(doc, vec))
95         return vec
96
97     def Compression(self, round_num=0) :
98         print(self.tdm)
99         print(self.term_doc_mat.round(round_num))
100
101     def TermVectorMatrix(self, u, s, topic_num):
102         term_mat = np.matrix(u[:, :topic_num])# * np.diag(s[:topic_num])
103         return term_mat
104
105     def DocVectorMatrix(self, s, vt, topic_num):
106         doc_mat = np.diag(s[:topic_num]) * np.matrix(vt[:topic_num,:])
107         return doc_mat
108
109     def TermDocVectorMatrix(self, u, s, vt, topic_num):

```

```

110 term_doc_mat = np.matrix(u[:, :topic_num]) * np.diag(s[:topic_num]) * np.matrix(vt[:topic_num,:])
111 return term_doc_mat
112
113 def TermSimilarityMatrix(self, termVectorMatrix):
114     return self.vectorSimilarity(termVectorMatrix.T)
115
116 def GetTermSimilarity(self, term1, term2):
117     sim = self.term_sim[self.term2idx[term1], self.term2idx[term2]]
118     print("{}{} term similarity = {}".format(term1, term2, sim))
119     return sim
120
121 def DocSimilarityMartrix(self, docVectorMatrix):
122     return self.vectorSimilarity(docVectorMatrix)
123
124 def GetDocSimilarity(self, doc1, doc2):
125     sim = self.doc_sim[self.doc2idx[doc1], self.doc2idx[doc2]]
126     print("{}{} doc similarity = {}".format(doc1, doc2, sim))
127     return sim

```

```

1 doc_ls = ['바나나 사과 포도 포도 짜장면',
2           '사과 포도',
3           '포도 바나나',
4           '짜장면 짬뽕 탕수육',
5           '볶음밥 탕수육',
6           '짜장면 짬뽕',
7           '라면 스시',
8           '스시 짜장면',
9           '가츠동 스시 소바',
10          '된장찌개 김치찌개 김치',
11          '김치 된장 짜장면',
12          '비빔밥 김치'
13          ]

```

```

14
15 lsa = LSA(doc_ls, 4)
16 X = lsa.TDM(doc_ls)
17 print('== 토픽 모델링 ==')
18 lsa.TopicModeling(4)
19 print('\n== 단어 벡터 ==')
20 lsa.GetTermVector('사과')
21 lsa.GetTermVector('짜장면')
22 print('\n== 단어 유사도 ==')
23 lsa.GetTermSimilarity('사과', '바나나')
24 lsa.GetTermSimilarity('사과', '짜장면')
25 lsa.GetTermSimilarity('포도', '짜장면')
26 lsa.GetTermSimilarity('사과', '스시')
27 print('\n== 문서 벡터 ==')
28 lsa.GetDocVector('사과 포도')
29 lsa.GetDocVector('짜장면 짬뽕')
30 print('\n== 문서 유사도 ==')

```

```
31 | Isa.GetDocSimilarity('사과 포도', '포도 바나나')
32 | Isa.GetDocSimilarity('사과 포도', '라면 스시')
33 | print('Wn== 토픽 차원수로 압축 ==')
34 | Isa.Compression(0)
```



Topic 1 - [('포도', 0.697), ('짜장면', 0.486), ('사과', 0.348), ('바나나', 0.348)]
 Topic 2 - [('짜장면', 0.584), ('짬뽕', 0.356), ('김치', 0.337), ('스시', 0.256)]
 Topic 3 - [('김치', 0.611), ('된장찌개', 0.264), ('김치찌개', 0.264), ('비빔밥', 0.185)]
 Topic 4 - [('스시', 0.552), ('김치', 0.371), ('가츠동', 0.277), ('소바', 0.277)]

사과 = [[0.34843127 -0.19370961 0.01592593 0.03744775]]
 짜장면 = [[0.48563449 0.58415588 -0.07468389 -0.18737521]]

(사과,바나나) term similarity = 1.0
(사과,짜장면) term similarity = 0.1519133521480699
(포도,짜장면) term similarity = 0.15191335214807006
(사과,스시) term similarity = -0.04097825202732752

사과 포도 = [[1.04529381 -0.58112882 0.04777778 0.11234324]]
 짜장면 짬뽕 = [[0.61011838 0.93971158 -0.17760018 -0.53682795]]

```
('사과 포도', '포도 바나나') doc similarity = 0.9999999999999998
('사과 포도', '라면 스시') doc similarity = -0.03850688211350069
```

```
[[1. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
 [1. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
 [2. 1. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
 [1. 0. 0. 1. 0. 1. 0. 1. 0. 0. 1. 0.]
 [0. 0. 0. 1. 0. 1. 0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 1. 1. 0. 0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 0. 1. 1. 1. 0. 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. 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. 0. 0. 1. 1. 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. 0. 0. 0. 1.]]
[[ 1.  0.  0. -0. -0.  0. -0.  0. -0.  0. -0.]
 [ 1.  0.  0. -0. -0.  0. -0.  0. -0. -0.  0. -0.]
 [ 2.  1.  1. -0. -0.  0. -0.  0. -0. -0.  0. -0.]
 [ 1.  0.  0.  1.  0.  1.  0.  1.  0.  0.  1.  0.]
 [ 0. -0. -0.  1.  0.  1. -0.  0. -0. -0.  0. -0.]
 [ 0. -0. -0.  1.  0.  0. -0.  0. -0. -0.  0. -0.]
 [-0. -0. -0.  0.  0.  0. -0.  0. -0. -0.  0. -0.]
 [-0. -0. -0. -0. -0. -0.  0.  0.  0. -0. -0. -0.]
 [ 0. -0. -0.  0. -0.  0.  1.  1.  1. -0.  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. -0. -0. -0.  0.  0.  0.]
 [-0. -0. -0. -0. -0. -0. -0. -0. -0.  0.  0.  0.]
 [ 0. -0. -0. -0. -0.  0. -0.  0. -0.  1.  1.  1.]
 [ 0. -0. -0.  0.  0.  0. -0.  0. -0.  0.  0.  0.]
 [-0. -0. -0. -0. -0. -0. -0. -0. -0.  0.  0.  0.]]
```

▼ 1.3 sklearn LSA 구현

```
1 doc_ls = ['바나나 사과 포도 포도',
2           '사과 포도',
3           '포도 바나나',
4           '짜장면 짬뽕 탕수육',
5           '볶음밥 탕수육',
6           '짜장면 짬뽕',
7           '라면 스시',
8           '스시',
9           '가츠동 스시 소바',
10          '된장찌개 김치찌개 김치',
11          '김치 된장',
12          '비빔밥 김치'
13          ]
```

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2
3 vectorizer = TfidfVectorizer(
4 max_features= 1000, # 상위 1,000개의 단어를 보존
5 max_df = 0.5,
```



```

6 smooth_idf=True)
7
8 X = vectorizer.fit_transform(doc_ls)
9
10 from sklearn.decomposition import TruncatedSVD
11 svd_model = TruncatedSVD(n_components=4, algorithm='randomized', n_iter=100)
12 svd_model.fit(X)
13
14 np.shape(svd_model.components_)
15
16 terms = vectorizer.get_feature_names() # 단어 집합. 1,000개의 단어가 저장됨.
17
18 def get_topics(components, feature_names, n=3):
19     for idx, topic in enumerate(components):
20         print("Topic %d:" % (idx+1), [(feature_names[i], topic[i].round(5)) for i in topic.argsort()[::-n - 1:-1]])
21 get_topics(svd_model.components_, terms)

```

Topic 1: [('포도', 0.78069), ('사과', 0.44189), ('바나나', 0.44189)]
 Topic 2: [('스시', 0.8864), ('라면', 0.33189), ('가츠동', 0.2282)]
 Topic 3: [('짬뽕', 0.6258), ('짜장면', 0.6258), ('탕수육', 0.43614)]
 Topic 4: [('김치', 0.76421), ('된장', 0.37119), ('비빔밥', 0.37119)]

2 LDA 실습

2.1 실습 템플릿

```

1 class LDA :
2     def __init__(self, doc_ls, topic_num, alpha = 0.1, beta = 0.001):
3         self.alpha = alpha
4         self.beta = beta
5         self.k = topic_num
6
7     def RandomlyAssignTopic(self, doc_ls):
8         pass
9
10
11     def IterateAssignTopic(self) :
12         pass
13
14
15     # 토픽별 주요 키워드 출력
16     def TopicModeling(self) :

```

```
17 | pass
18 |
```

▼ 2.2 실습 예제 코드

```
1 | import nltk
2 | nltk.download('punkt')
3 | nltk.download('stopwords')
4 | nltk.download('wordnet')
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Unzipping corpora/wordnet.zip.
True
```

```
1 | import random
2 | import numpy as np
3 | from nltk.tokenize import word_tokenize
4 | from nltk.stem import WordNetLemmatizer
5 | from nltk.corpus import stopwords
6 | from collections import defaultdict
7 |
8 |
9 |
10 | class LDA :
11 |     def __init__(self, docs, topic_num, alpha = 0.1, beta = 0.001):
12 |         self.alpha = alpha
13 |         self.beta = beta
14 |         self.k = topic_num
15 |         self.docs = docs
16 |
17 |
18 |     def RandomlyAssignTopic(self, docs) :
19 |         dic = defaultdict(dict)
20 |         t2i = defaultdict(lambda : len(t2i))
21 |         i2t = defaultdict()
22 |         d = 0
23 |         w = 0
24 |
25 |         wnl = WordNetLemmatizer()
```

```

26  stopwords = stopwords.words('english')
27  stopwords.append(',')
28
29  # 임의의 토픽을 할당
30  for tokens in [word_tokenize(doc) for doc in docs] :
31      for token in [wnl.lemmatize(token.lower()) for token in tokens if token not in stopwords] :
32          i2t[t2i[token]] = token
33          dic[(d, t2i[token], w)] = random.randint(0, self.k-1)
34          w += 1
35          d += 1
36
37  #print(dic)
38  return dic, t2i, i2t
39
40
41  def CountDocTopic(self) :
42      docs = np.zeros((self.k, len(self.docs)))
43      terms = np.zeros((self.k, len(self.t2i.keys())))
44
45
46      #문서별 토큰별 빈도수 계산
47      for (d, n, w) in self.term_topic.keys() :
48          docs[self.term_topic[(d, n, w)], d] += 1 + self.alpha
49          terms[self.term_topic[(d, n, w)], d] += 1 + self.beta
50      #print(doc_m)
51
52      #비어있는 값은 값에 alpha, beta 설정
53      docs = np.where(docs==0.0, self.alpha, docs)
54      terms = np.where(terms==0.0, self.beta, terms)
55
56      return docs, terms
57
58
59  def IterateAssignTopic(self, docs, terms) :
60      #한개의 단어씩 주제 배정
61      prev = {}
62
63      while prev != self.term_topic:
64          for (d, n, w) in self.term_topic.keys() :
65              topic = [0, 0]
66
67              docs[self.term_topic[(d, n, w)], d] -= (1 + self.alpha)
68              terms[self.term_topic[(d, n, w)], n] -= (1 + self.beta)
69
70              #print(doc_m)
71              prev = self.term_topic
72
73              for t in range(self.k) :
74                  p_t_d = docs[t, d]/docs[:,d:d+1].sum()
75                  p_w_t = terms[t, n]/terms[t:t+1,:].sum()
76                  prob = p_t_d * p_w_t

```

```

77
78         if topic[1] < prob :
79             topic = [t, prob]
80
81         #확률이 가장 큰 토픽을 할당
82         self.term_topic[(d, n, w)] = topic[0]
83         docs[self.term_topic[(d, n, w)], d] += (1 + self.alpha)
84         terms[self.term_topic[(d, n, w)], t] += (1 + self.beta)
85
86         #print(self.term_topic)
87
88     return terms
89
90
91 # 토픽별 주요 키워드 출력
92 def TopicModeling(self, count=5) :
93     self.term_topic, self.t2i, self.i2t = self.RandomlyAssignTopic(self.docs)
94     docs, terms = self.CountDocTopic()
95     terms = self.IterateAssignTopic(docs, terms)
96
97     score = terms / terms.sum(axis=1, keepdims=True)
98
99     for i in range(self.k) :
100         print("\nTopic {}".format(i+1))
101         sorted_index = np.flip(np.argsort(score[i]),0)[:count]
102         for j in sorted_index :
103             print("{}={}".format(self.i2t[j], score[i,j].round(3)), end = ' ')
104
105

```

```

1 doc_ls = ["Cute kitty",
2 "Eat rice or cake",
3 "Kitty and hamster",
4 "Eat bread",
5 "Rice, bread and cake",
6 "Cute hamster eats bread and cake"]
7
8 lda = LDA(doc_ls, 2)
9 lda.TopicModeling()

```



▼ 2.3 pyLDAvis 를 이용한 LDA 시각화

```
1 !pip install pyLDAvis
```



```

1 import pandas as pd
2 from sklearn.datasets import fetch_20newsgroups
3 dataset = fetch_20newsgroups(shuffle=True, random_state=1, remove=('headers', 'footers', 'quotes'))
4 documents = dataset.data
5 len(documents)

```



```

1 news_df = pd.DataFrame({'document':documents})
2 # 특수 문자 제거
3 news_df['clean_doc'] = news_df['document'].str.replace("[^a-zA-Z]", " ")
4 # 길이가 3이하인 단어는 제거 (길이가 짧은 단어 제거)
5 news_df['clean_doc'] = news_df['clean_doc'].apply(lambda x: ' '.join([w for w in x.split() if len(w)>3]))
6 # 전체 단어에 대한 소문자 변환
7 news_df['clean_doc'] = news_df['clean_doc'].apply(lambda x: x.lower())

```

```

1 import nltk
2 nltk.download('stopwords')
3 from nltk.corpus import stopwords
4 stop_words = stopwords.words('english') # NLTK로부터 불용어를 받아옵니다.
5 tokenized_doc = news_df['clean_doc'].apply(lambda x: x.split()) # 토큰화
6 tokenized_doc = tokenized_doc.apply(lambda x: [item for item in x if item not in stop_words])

```



```

1 from gensim import corpora
2 dictionary = corpora.Dictionary(tokenized_doc)
3 corpus = [dictionary.doc2bow(text) for text in tokenized_doc]
4 print(corpus[1]) # 수행된 결과에서 두번째 뉴스 출력. 첫번째 문서의 인덱스는 0

```



```

1 import gensim
2 NUM_TOPICS = 20 #20개의 토픽, k=20
3 ldamodel = gensim.models.ldamodel.LdaModel(corpus, num_topics = NUM_TOPICS, id2word=dictionary, passes=15)
4 topics = ldamodel.print_topics(num_words=4)
5 for topic in topics:

```

```
6 | print(topic)
```

```
|
```



```
1 | import pyLDAvis.gensim
2 | pyLDAvis.enable_notebook()
3 | vis = pyLDAvis.gensim.prepare(ldamodel, corpus, dictionary)
4 | pyLDAvis.display(vis)
```



```

1 for i, topic_list in enumerate(ldamodel[corpus]):
2     if i==5:
3         break
4     print(i, '번째 문서의 topic 비율은', topic_list)

1 def make_topictable_per_doc(ldamodel, corpus, texts):
2     topic_table = pd.DataFrame()
3
4     # 몇 번째 문서인지를 의미하는 문서 번호와 해당 문서의 토픽 비중을 한 줄씩 꺼내온다.
5     for i, topic_list in enumerate(ldamodel[corpus]):
6         doc = topic_list[0] if ldamodel.per_word_topics else topic_list
7         doc = sorted(doc, key=lambda x: (x[1]), reverse=True)
8         # 각 문서에 대해서 비중이 높은 토픽순으로 토픽을 정렬한다.
9         # EX) 정렬 전 0번 문서 : (2번 토픽, 48.5%), (8번 토픽, 25%), (10번 토픽, 5%), (12번 토픽, 21.5%),
10        # EX) 정렬 후 0번 문서 : (2번 토픽, 48.5%), (8번 토픽, 25%), (12번 토픽, 21.5%), (10번 토픽, 5%)
11        # 48 > 25 > 21 > 5 순으로 정렬이 된 것.
12
13        # 모든 문서에 대해서 각각 아래를 수행
14        for j, (topic_num, prop_topic) in enumerate(doc): # 몇 번 토픽인지와 비중을 나눠서 저장한다.
15            if j == 0: # 정렬을 한 상태이므로 가장 앞에 있는 것이 가장 비중이 높은 토픽
16                topic_table = topic_table.append(pd.Series([int(topic_num), round(prop_topic, 4), topic_list]), ignore_index=True)
17                # 가장 비중이 높은 토픽과, 가장 비중이 높은 토픽의 비중과, 전체 토픽의 비중을 저장한다.
18            else:
19                break
20    return(topic_table)

1 topictable = make_topictable_per_doc(ldamodel, corpus, tokenized_doc)
2 topictable = topictable.reset_index() # 문서 번호를 의미하는 열(column)로 사용하기 위해서 인덱스 열을 하나 더 만든다.
3 topictable.columns = ['문서 번호', '가장 비중이 높은 토픽', '가장 높은 토픽의 비중', '각 토픽의 비중']
4 topictable[:10]

```



| 문서 번호 | 가장 비중이 높은 토픽 | 가장 높은 토픽의 비중 | 각 토픽의 비중 |
|----------|--------------|--------------|--|
| 0 | 0 | 17.0 | 0.6845 [(8, 0.30094197), (17, 0.6845419)] |
| 1 | 1 | 8.0 | 0.5632 [(2, 0.029335152), (8, 0.5632338), (10, 0.1083... |
| 2 | 2 | 8.0 | 0.6275 [(4, 0.018661937), (8, 0.62747276), (16, 0.045... |
| 3 | 3 | 19.0 | 0.2927 [(1, 0.044962943), (8, 0.14445594), (9, 0.0294... |
| 4 | 4 | 10.0 | 0.6536 [(2, 0.082709715), (5, 0.23224762), (10, 0.653... |
| 5 | 5 | 8.0 | 0.6517 [(0, 0.105608456), (3, 0.15081134), (8, 0.6517... |
| 6 | 6 | 17.0 | 0.7437 [(2, 0.013339129), (7, 0.012894245), (8, 0.014... |
| 7 | 7 | 8.0 | 0.4894 [(8, 0.48941994), (10, 0.12249057), (13, 0.024... |
| 8 | 8 | 3.0 | 0.3653 [(0, 0.2194796), (3, 0.36532196), (8, 0.220569... |
| 9 | 9 | 19.0 | 0.4796 [(0, 0.038544882), (4, 0.04621158), (6, 0.0463... |

