

▼ TF-IDF

▼ scikit-learn을 활용한 TF-IDF구현

```
1 from sklearn.feature_extraction.text import CountVectorizer
```

```
1 corpus = [
2     'you know I want your love',
3     'I like you',
4     'what should I do'
5 ]
```

```
1 vector = CountVectorizer()
```

```
1 print(vector.fit_transform(corpus).toarray())
```

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

```
1 print(vector.vocabulary_)
```

```
{'you': 7, 'know': 1, 'want': 5, 'your': 8, 'love': 3, 'like': 2, 'what': 6, 'I': 0, 'should': 4}
```

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
```

```
1 corpus = [
2     'you know I want your love',
3     'I like you',
4     'what should I do'
5 ]
```

```
1 tfidf = TfidfVectorizer().fit(corpus)
```

```
1 print(tfidf.transform(corpus).toarray())
```

```
[[0.          0.46735098 0.          0.46735098 0.          0.46735098
  0.          0.35543247 0.46735098]
 [0.          0.          0.79596054 0.          0.          0.
  0.          0.60534851 0.          ]
 [0.57735027 0.          0.          0.          0.57735027 0.
  0.57735027 0.          0.          ]]
```

```
1 print(tfidf.vocabulary_)
```

```
{'you': 7, 'know': 1, 'want': 5, 'your': 8, 'love': 3, 'like': 2, 'what': 6, '}
```

▼ 로이터 데이터로 TF-IDF 학습하기

```
1 from tensorflow.keras.datasets import reuters
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5 import pandas as pd
```

```
1 (x_train, y_train), (x_test, y_test) = reuters.load_data(num_words=10000, test_s
```

```
    Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datas
    2113536/2110848 [=====] - 0s 0us/step
    2121728/2110848 [=====] - 0s 0us/step
```

```
1 print('훈련 샘플의 수 : {}'.format(len(x_train)))
```

```
    훈련 샘플의 수 : 8982
```

```
1 print('테스트 샘플의 수 : {}'.format(len(x_test)))
```

```
    테스트 샘플의 수 : 2246
```

```
1 print(x_train[0])
2 print(x_test[0])
```

```
    [1, 2, 2, 8, 43, 10, 447, 5, 25, 207, 270, 5, 3095, 111, 16, 369, 186, 90, 67,
    [1, 4, 1378, 2025, 9, 697, 4622, 111, 8, 25, 109, 29, 3650, 11, 150, 244, 364,
```

```
1 print(y_train[0]) # 훈련 기사의 레이블
```

```
    3
```

```
1 num_classes = max(y_train)+1
2 print('클래스의 수 : {}'.format(num_classes))
```

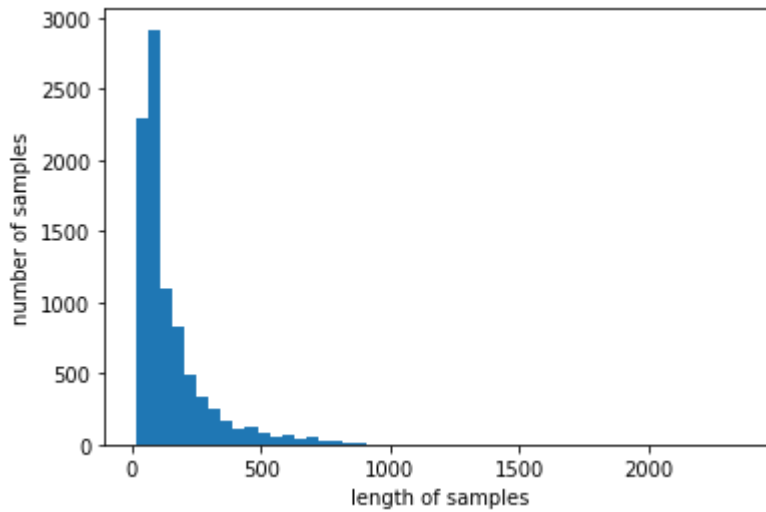
```
    클래스의 수 : 46
```

```
1 print('훈련용 뉴스의 최대 길이 : {}'.format(max(len(l) for l in x_train)))
2 print('훈련용 뉴스의 평균 길이 : {}'.format(sum(map(len, x_train))/len(x_train)))
```

```
    훈련용 뉴스의 최대 길이 : 2376
    훈련용 뉴스의 평균 길이 : 145.5398574927633
```

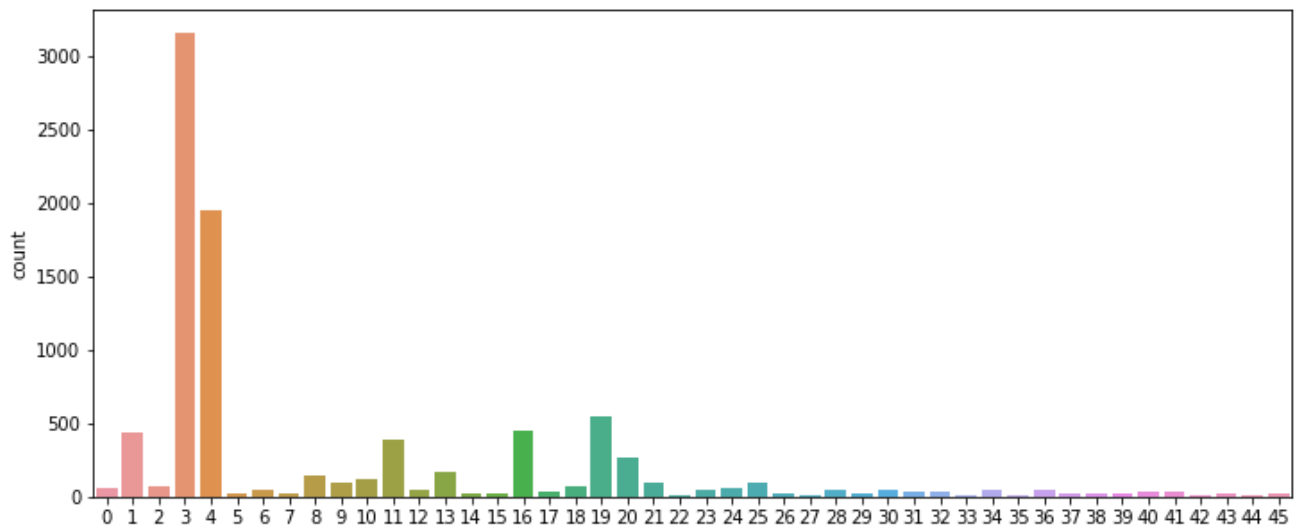
```
1 plt.hist([len(s) for s in x_train], bins= 50)
2 plt.xlabel('length of samples')
3 plt.ylabel('number of samples')
```

```
4 plt.show()
```



```
1 fig, axe = plt.subplots(ncols =1)
2 fig.set_size_inches(12, 5)
3 sns.countplot(y_train)
```

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning
FutureWarning
<matplotlib.axes._subplots.AxesSubplot at 0x7fe8a2d8f310>



```
1 unique_elements, counts_elements = np.unique(y_train, return_counts=True)
2 print("각 클래스 빈도수 : ")
3 print(np.asarray((unique_elements, counts_elements)))
```

각 클래스 빈도수 :

```
[[ 0  1  2  3  4  5  6  7  8  9 10 11 12 13
 14 15 16 17 18 19 20 21 22 23 24 25 26 27
 28 29 30 31 32 33 34 35 36 37 38 39 40 41
 42 43 44 45]
 [ 55 432  74 3159 1949  17  48  16 139 101 124 390  49 172
 26  20 444  39  66 549 269 100  15  41  62  92  24  15
 48  19  45  39  32  11  50  10  49  19  19  24  36  30
 13  21  12  18]]
```

```
1 word_to_index = reuters.get_word_index()
2 print(word_to_index)
```

```
Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datas
557056/550378 [=====] - 0s 0us/step
565248/550378 [=====] - 0s 0us/step
{'mdb1': 10996, 'fawc': 16260, 'degussa': 12089, 'woods': 8803, 'hanging': 137
```

```
1 index_to_word = {index + 3 : word for word, index in word_to_index.items()}
```

```
1 word_to_index['the']
```

```
1
```

```
1 word_to_index['it']
```

```
13
```

```
1 print(index_to_word[4]) # 빈도수 상위 1번 단어
2 print(index_to_word[16]) # 빈도수 상위 13번 단어
```

```
the
it
```

```
1 # 0 <pad>
2 # 1 <sos>
3 # 2 <unk>
4 for index, token in enumerate("<pad>", "<sos>", "<unk>"):
5     index_to_word[index] = token
```

```
1 print(' '.join([index_to_word[index] for index in x_train[0]]))
```

```
<sos> <unk> <unk> said as a result of its december acquisition of space co it
```

```
1 # 전체 훈련 데이터에 대해서 decoded
2 decoded = []
3 for i in range(len(x_train)):
4     t = ' '.join([index_to_word[index] for index in x_train[i]])
5     decoded.append(t)
6
7 x_train = decoded
```

```
1 # 전체 테스트 데이터에 대해서 decoded
2 decoded = []
3 for i in range(len(x_test)):
4     t = ' '.join([index_to_word[index] for index in x_test[i]])
5     decoded.append(t)
6
7 x_test = decoded
```

```
1 x_train[:5]
```

```
['<sos> <unk> <unk> said as a result of its december acquisition of space co i
<sos> generale de banque sa lt <unk> br and lt heller overseas corp of chica
<sos> shr 3 28 dlrs vs 22 cts shr diluted 2 99 dlrs vs 22 cts net 46 0 mln v
"<sos> the farmers home administration the u s agriculture department's farm
<sos> seton co said its board has received a proposal from chairman and chie
```

```
1 x_test[:5]
```

```
['<sos> the great atlantic and pacific tea co said its three year 345 mln dlr
"<sos> philippine sugar production in the 1987 88 crop year ending august has
"<sos> the agriculture department's widening of louisiana gulf differentials
<sos> <unk> <unk> oil and gas partnership said it completed the sale of inte
<sos> strong south <unk> winds were keeping many vessels trapped in the ice
```

▼ TF-IDF로 로이터 데이터 학습하기

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.feature_extraction.text import CountVectorizer
3 from sklearn.feature_extraction.text import TfidfTransformer
4 from sklearn.naive_bayes import MultinomialNB
5 from sklearn import metrics
```

```
1 # 단어의 수를 카운트하는 사이킷런의 카운트벡터라이저.
2 count_vect = CountVectorizer()
3 # fit_transform : 학습 할 때와 동일한 기반 설정으로 동일하게 테스트 데이터를 변환해야 하는 것
4 x_train_counts = count_vect.fit_transform(x_train)
5
6 # 카운트벡터라이저의 결과로부터 TF-IDF 결과를 얻습니다.
7 tfidf_transformer = TfidfTransformer()
8 x_train_tfidf = tfidf_transformer.fit_transform(x_train_counts)
9
10 # 나이브 베이즈 분류기를 수행
11 # x_train은 TF-IDF의 벡터, y_train 레이블
12 clf = MultinomialNB().fit(x_train_tfidf, y_train)
```

```
1 def tfidf_vectorizer(data):
2     data_counts = count_vect.transform(data)
3     data_tfidf = tfidf_transformer.transform(data_counts)
4     return data_tfidf
```

```
1 y_pred = clf.predict(tfidf_vectorizer(x_test))
```

```
1 print(metrics.classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.00	0.00	0.00	12

1	0.62	0.69	0.65	105			
2	0.00	0.00	0.00	20			
3	0.81	0.90	0.85	813			
4	0.51	0.96	0.67	474			
5	0.00	0.00	0.00	5			
6	0.00	0.00	0.00	14			
7	0.00	0.00	0.00	3			
8	0.00	0.00	0.00	38			
9	1.00	0.08	0.15	25			
10	0.00	0.00	0.00	30			
11	0.66	0.63	0.64	83			
12	0.00	0.00	0.00	13			
13	1.00	0.03	0.05	37			
14	0.00	0.00	0.00	2			
15	0.00	0.00	0.00	9			
16	0.69	0.56	0.61	99			
17	0.00	0.00	0.00	12			
18	0.00	0.00	0.00	20			
19	0.60	0.78	0.68	133			
20	1.00	0.04	0.08	70			
21	0.00	0.00	0.00	27			
22	0.00	0.00	0.00	7			
23	0.00	0.00	0.00	12			
24	0.00	0.00	0.00	19			
25	1.00	0.03	0.06	31			
26	0.00	0.00	0.00	8			
27	0.00	0.00	0.00	4			
28	0.00	0.00	0.00	10			
29	0.00	0.00	0.00	4			
30	0.00	0.00	0.00	12			
31	0.00	0.00	0.00	13			
32	0.00	0.00	0.00	10			
33	0.00	0.00	0.00	5			
34	0.00	0.00	0.00	7			
35	0.00	0.00	0.00	6			
36	0.00	0.00	0.00	11			
37	0.00	0.00	0.00	2			
38	0.00	0.00	0.00	3			
39	0.00	0.00	0.00	5			
40	0.00	0.00	0.00	10			
41	0.00	0.00	0.00	8			
42	0.00	0.00	0.00	3			
43	0.00	0.00	0.00	6			
44	0.00	0.00	0.00	5			
45	0.00	0.00	0.00	1			
accuracy				0.66	2246		
macro avg				0.17	0.10	0.10	2246
weighted avg				0.59	0.66	0.58	2246

```

/usr/local/lib/python3.7/dist-packages/sklearn/metrics/_classification.py:1308
_warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.7/dist-packages/sklearn/metrics/_classification.py:1308
_warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.7/dist-packages/sklearn/metrics/_classification.py:1308

```

Lstm모델 구현

```
1 from tensorflow.keras.models import Sequential
```

```

2 from tensorflow.keras.layers import Dense, LSTM, Embedding
3 from tensorflow.keras.preprocessing.sequence import pad_sequences
4 from tensorflow.keras.utils import to_categorical
5 from tensorflow.keras.callbacks import EarlyStopping, ModelCheckpoint
6 from tensorflow.keras.models import load_model

1 (x_train, y_train), (x_test, y_test) = reuters.load_data(num_words=1000, test_sp

1 max_len = 100
2 x_train = pad_sequences(x_train, maxlen= max_len)
3 x_test = pad_sequences(x_test, maxlen=max_len)

1 y_train = to_categorical(y_train)
2 y_test = to_categorical(y_test)

1 vocab_size = 1000
2 embedding_dim = 128
3 hidden_units = 128
4 num_classes = 46
5
6 model = Sequential()
7 model.add(Embedding(vocab_size, embedding_dim))
8 model.add(LSTM(hidden_units))
9 model.add(Dense(num_classes, activation='softmax'))

1 es = EarlyStopping(monitor='val_loss', mode='min', verbose=1, patience= 4)
2 mc = ModelCheckpoint('best_model.h5', monitor='val_acc', mode='max', verbose = 1

1 model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['acc'])

1 history = model.fit(x_train, y_train, batch_size = 128, epochs= 30, callbacks=[e

Epoch 1/30
71/71 [=====] - ETA: 0s - loss: 2.6258 - acc: 0.3477
Epoch 00001: val_acc improved from -inf to 0.36198, saving model to best_model
71/71 [=====] - 21s 274ms/step - loss: 2.6258 - acc:
Epoch 2/30
71/71 [=====] - ETA: 0s - loss: 2.1235 - acc: 0.4591
Epoch 00002: val_acc improved from 0.36198 to 0.50000, saving model to best_mc
71/71 [=====] - 19s 268ms/step - loss: 2.1235 - acc:
Epoch 3/30
71/71 [=====] - ETA: 0s - loss: 1.8979 - acc: 0.5125
Epoch 00003: val_acc improved from 0.50000 to 0.53117, saving model to best_mc
71/71 [=====] - 19s 268ms/step - loss: 1.8979 - acc:
Epoch 4/30
71/71 [=====] - ETA: 0s - loss: 1.7365 - acc: 0.5530
Epoch 00004: val_acc improved from 0.53117 to 0.56055, saving model to best_mc
71/71 [=====] - 19s 268ms/step - loss: 1.7365 - acc:
Epoch 5/30
71/71 [=====] - ETA: 0s - loss: 1.6632 - acc: 0.5757
Epoch 00005: val_acc did not improve from 0.56055
71/71 [=====] - 19s 269ms/step - loss: 1.6632 - acc:

```

```

Epoch 6/30
71/71 [=====] - ETA: 0s - loss: 1.6406 - acc: 0.5835
Epoch 00006: val_acc improved from 0.56055 to 0.56679, saving model to best_model.h5
71/71 [=====] - 19s 269ms/step - loss: 1.6406 - acc: 0.5835
Epoch 7/30
71/71 [=====] - ETA: 0s - loss: 1.5375 - acc: 0.6100
Epoch 00007: val_acc improved from 0.56679 to 0.60686, saving model to best_model.h5
71/71 [=====] - 19s 269ms/step - loss: 1.5375 - acc: 0.6100
Epoch 8/30
71/71 [=====] - ETA: 0s - loss: 1.4402 - acc: 0.6342
Epoch 00008: val_acc improved from 0.60686 to 0.62556, saving model to best_model.h5
71/71 [=====] - 19s 268ms/step - loss: 1.4402 - acc: 0.6342
Epoch 9/30
71/71 [=====] - ETA: 0s - loss: 1.3486 - acc: 0.6548
Epoch 00009: val_acc improved from 0.62556 to 0.63090, saving model to best_model.h5
71/71 [=====] - 19s 270ms/step - loss: 1.3486 - acc: 0.6548
Epoch 10/30
71/71 [=====] - ETA: 0s - loss: 1.2875 - acc: 0.6740
Epoch 00010: val_acc improved from 0.63090 to 0.64871, saving model to best_model.h5
71/71 [=====] - 19s 268ms/step - loss: 1.2875 - acc: 0.6740
Epoch 11/30
71/71 [=====] - ETA: 0s - loss: 1.2007 - acc: 0.6981
Epoch 00011: val_acc improved from 0.64871 to 0.66429, saving model to best_model.h5
71/71 [=====] - 19s 269ms/step - loss: 1.2007 - acc: 0.6981
Epoch 12/30
71/71 [=====] - ETA: 0s - loss: 1.1551 - acc: 0.7050
Epoch 00012: val_acc improved from 0.66429 to 0.67097, saving model to best_model.h5
71/71 [=====] - 19s 269ms/step - loss: 1.1551 - acc: 0.7050
Epoch 13/30
71/71 [=====] - ETA: 0s - loss: 1.0824 - acc: 0.7244
Epoch 00013: val_acc improved from 0.67097 to 0.67231, saving model to best_model.h5
71/71 [=====] - 19s 269ms/step - loss: 1.0824 - acc: 0.7244
Epoch 14/30
71/71 [=====] - ETA: 0s - loss: 1.0375 - acc: 0.7337
Epoch 00014: val_acc improved from 0.67231 to 0.67409, saving model to best_model.h5
71/71 [=====] - 19s 268ms/step - loss: 1.0375 - acc: 0.7337
Epoch 15/30
71/71 [=====] - ETA: 0s - loss: 0.9899 - acc: 0.7518

```

```

1 loaded_model = load_model('best_model.h5')
2 print("\n 테스트 정확도 : %.4f" % (loaded_model.evaluate(x_test, y_test)[1]))

71/71 [=====] - 2s 24ms/step - loss: 1.2067 - acc: 0.7168

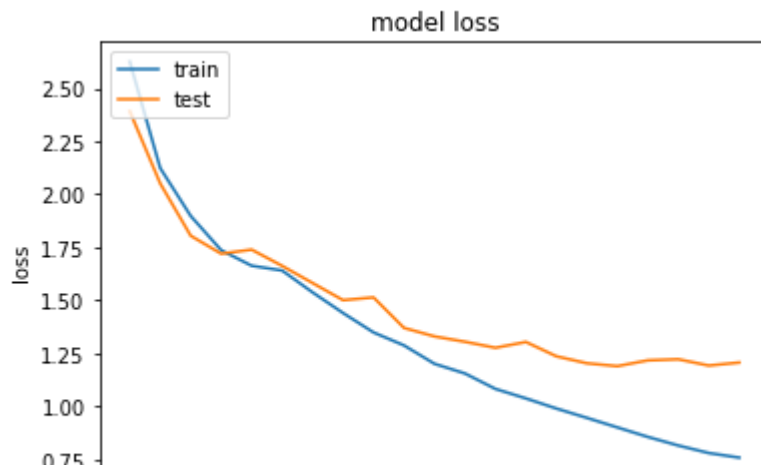
테스트 정확도 : 0.7168

```

```

1 epochs = range(1, len(history.history['acc'])+1)
2 plt.plot(epochs, history.history['loss'])
3 plt.plot(epochs, history.history['val_loss'])
4 plt.title('model loss')
5 plt.ylabel('loss')
6 plt.xlabel('epoch')
7 plt.legend(['train', 'test'], loc='upper left')
8 plt.show()

```

▼ 실제 뉴스 기사 크롤링 및 분류

```
1 !pip install beautifulsoup4
2 !pip install newspaper3k
3 !pip install konlpy
```

[illegible]


```

9         <span class = 'price'> 29000 </span>
10        <span class = 'menu'> 의류 </span>
11        <a href = "http://www.naver.com"> 바로가기 </a>
12    </p>
13    <p id='watch' class= 'name' title = '시계'> 시계
14        <span class = 'number'> 28 </span>
15        <span class = 'price'> 32000 </span>
16        <span class = 'menu'> 악세사리 </span>
17        <a href = "http://www.facebook.com"> 바로가기 </a>
18    </p>
19 </h1>
20 </body>
21 </html>
22 '''

```

```
1 soup = BeautifulSoup(html, 'html.parser')
```

```
1 print(soup.select('body'))
```

```

[<body>
<h1>장바구니
    <p class="name" id="clothes" title="라운드티"> 라운드티
        <span class="number"> 25 </span>
<span class="price"> 29000 </span>
<span class="menu"> 의류 </span>
<a href="http://www.naver.com"> 바로가기 </a>
</p>
<p class="name" id="watch" title="시계"> 시계
    <span class="number"> 28 </span>
<span class="price"> 32000 </span>
<span class="menu"> 악세사리 </span>
<a href="http://www.facebook.com"> 바로가기 </a>
</p>
</h1>
</body>]

```

```
1 print(soup.select('p'))
```

```

[<p class="name" id="clothes" title="라운드티"> 라운드티
    <span class="number"> 25 </span>
<span class="price"> 29000 </span>
<span class="menu"> 의류 </span>
<a href="http://www.naver.com"> 바로가기 </a>
</p>, <p class="name" id="watch" title="시계"> 시계
    <span class="number"> 28 </span>
<span class="price"> 32000 </span>
<span class="menu"> 악세사리 </span>
<a href="http://www.facebook.com"> 바로가기 </a>
</p>]

```

```
1 print(soup.select('h1 .name .menu'))
```

```
[<span class="menu"> 의류 </span>, <span class="menu"> 악세사리 </span>]
```

```
1 print(soup.select('html > h1'))
```

```
1 print(soup.select('h1'))
```

```
[]
```

▼ Newspaper3k 패키지

```
1 from newspaper import Article
```

```
1 url = 'https://news.naver.com/main/read.naver?mode=LSD&mid=sec&sid1=101&oid=030&
```

```
1 article = Article(url, language='ko')
```

```
2 article.download()
```

```
3 article.parse()
```

```
1 print('기사 제목 : ')
```

```
2 print(article.title)
```

기사 제목 :

[AI 사피엔스 시대]자연어처리 기술, 컴퓨팅 파워 경쟁 시대로

```
1 print('기사 내용 :')
```

```
2 print(article.text)
```

기사 내용 :

[Copyright © 전자신문 & 전자신문인터넷, 무단전재 및 재배포 금지]

주로 아이디어와 기술력으로 경쟁했던 자연어처리 인공지능(AI) 분야는 점차 컴퓨팅 파워 싸움으로 무게 추가

▼ BeautifulSoup와 newspaper3k를 통해 크롤러 만들기

```
1 import requests
```

```
2 import pandas as pd
```

```
3 from bs4 import BeautifulSoup
```

```
1 def make_urllist(page_num, code, date):
```

```
2     urllist= []
```

```
3     for i in range(1, page_num + 1):
```

```
4         url = 'https://news.naver.com/main/list.nhn?mode=LSD&mid=sec&sid1='+str(code)
```

```
5         headers = {'User-Agent': 'Mozilla/5.0 (X11; Linux x86_64) AppleWebKit/537.36
```

```
6         news = requests.get(url, headers=headers)
```

```
7
```

```
8     soup = BeautifulSoup(news.content, 'html.parser')
```

```
9
```

```
10     # CASE 1
```

```
11     news_list = soup.select('.newsflash_body .type06_headline li dl')
```

```
12     # CASE 2
```

```
13     news_list.extend(soup.select('.newsflash_body .type06 li dl'))
```

```
14
```

```
15     # 각 뉴스로부터 태그의 href를 추출해서 리스트에 추가합니다
```

```

15     # 각 페이지에서 a 태그의 <a href = ... > 에는 url 링크가 담겨있다.
16     for line in news_list:
17         urllist.append(line.a.get('href'))
18     return urllist

```

```

1 url_list = make_urllist(2, 101, 20200506)
2 print('뉴스 기사의 갯수 :', len(url_list))

```

뉴스 기사의 갯수 : 40

```
1 url_list[:5]
```

```

[ 'https://news.naver.com/main/read.naver?mode=LSD&mid=sec&sid1=101&oid=057&aic
' https://news.naver.com/main/read.naver?mode=LSD&mid=sec&sid1=101&oid=057&aic
' https://news.naver.com/main/read.naver?mode=LSD&mid=sec&sid1=101&oid=057&aic
' https://news.naver.com/main/read.naver?mode=LSD&mid=sec&sid1=101&oid=003&aic
' https://news.naver.com/main/read.naver?mode=LSD&mid=sec&sid1=101&oid=057&aic

```

```
1 idx2word = {'101' : '경제', '102': '사회', '103' : '생활/문화', '105' : 'IT/과학'}
```

```

1 from newspaper import Article
2
3 # 데이터프레임을 생성하는 함수
4 def make_date(urllist, code):
5     text_list = []
6     for url in urllist:
7         article = Article(url, language='ko')
8         article.download()
9         article.parse()
10        text_list.append(article.text)
11
12    # 데이터프레임의 'news'티 아래 파싱한 텍스트를 밸류로 붙여준다.
13    df = pd.DataFrame({'news':text_list})
14
15    # 데이터프레임의 'code'키 아래 한글 카테고리명을 붙여준다.
16    df['code'] = idx2word[str(code)]
17    return df

```

```

1 data = make_date(url_list, 101)
2 data[:10]

```

	news	code
0	고려은단이 5월을 맞아 응원 메시지를 공유하는 '5월 5글자로 응원 부탁해!' 이벤...	경제
1	코리아나화장품의 민감성 피부를 위한 저자극 스킨케어 브랜드 '프리엔제'가 마르고 건...	경제
2	서울장수주식회사가 부드럽고 달콤한 맛으로 인기를 모으고 있는 생막걸리 '인생막걸리'...	경제

▼ 데이터 수집 및 전처리

```

1 code_list = [102, 103, 105]
2 code_list

[102, 103, 105]

1 def make_total_data(page_num, code_list, date):
2     df = None
3
4     for code in code_list:
5         url_list = make_urllist(page_num, code, date)
6         df_temp = make_date(url_list, code)
7         print(str(code) + '번 코드에 대한 데이터를 만들었습니다.')
8
9         if df is not None:
10             df = pd.concat([df, df_temp])
11         else:
12             df = df_temp
13
14     return df

1 df = make_total_data(1, code_list, 20200506)

102번 코드에 대한 데이터를 만들었습니다.
103번 코드에 대한 데이터를 만들었습니다.
105번 코드에 대한 데이터를 만들었습니다.

1 print('뉴스 기사의 갯수 :', len(df))

뉴스 기사의 갯수 : 60

1 df.sample(10)

```

	news	code
19	넷플릭스와 유튜브, 페이스북 등에게 국내 이용자를 위한 '서비스 안정성'을 유지할 ...	IT/과학
2	황범순 의정부시 부시장 을지대학교 의정부캠퍼스 및 부속병원 공사현장 안전점검. 사진...	사회
18	기사 섹션 분류 안내\n\n기사의 섹션 정보는 해당 언론사의 분류를 따르고 있습니다...	IT/과학
10	기사 섹션 분류 안내\n\n기사의 섹션 정보는 해당 언론사의 분류를 따르고 있습니다...	생활/문화
13	블랙홀을 품은 삼중성계 HR 6819 상상도 [ESO/L. Calçada 제공/ 재...	IT/과학

```

1 df = make_total_data(100, code_list, 20200506)

```

102번 코드에 대한 데이터를 만들었습니다.
 103번 코드에 대한 데이터를 만들었습니다.
 105번 코드에 대한 데이터를 만들었습니다.

	news	code
4	지난 2016년 포항공대에 구축한 4세대 선형 방사광가속기. /연합뉴스 지난 201...	IT/과학

```

1 import os
2
3 csv_path = '../news_data.csv'
4 df.to_csv(csv_path, index=False)

1 if os.path.exists(csv_path):
2     print('{} File Saved!'.format(csv_path))

../news_data.csv File Saved!

1 #csv_path = '구글드라이브 주소 경로'
2 df = pd.read_table(csv_path, sep=',')
3 df.head()

```

	news	code
0	파주시청. 사진제공=파주시 파주시청. 사진제공=파주시\n\n[파주=파이낸셜뉴스 강근...	사회
1	동영상 뉴스\n\n이천 물류창고 화재 발화지점으로 지목된 지하 2층에서 산소절단기의...	사회
2	황범순 의정부시 부시장 을지대학교 의정부캠퍼스 및 부속병원 공사현장 안전점검. 사진...	사회
3	귀갓길 여성을 쫓아가 성범죄를 시도한 20대 남성이 구속됐습니다.서울 강남경찰서는 ...	사회
4	(서울=연합뉴스) 대한약사회가 6일부터 코로나바이러스 감염증 대응 체계를 '사회적 ...	사회

```

1 df['news'] = df['news'].str.replace("[^ㄱ-ㅎㅏ-ㅣ가-힣 ]", "")
2 df['news']

0      파주시청사진제공파주시 파주시청사진제공파주시파주파이낸셜뉴스 강근주 기자파주시는 일 관...
1      동영상 뉴스이천 물류창고 화재 발화지점으로 지목된 지하 층에서 산소절단기의 산소 공...
2      황범순 의정부시 부시장 을지대학교 의정부캠퍼스 및 부속병원 공사현장 안전점검사진제공...
3      귀갓길 여성을 쫓아가 성범죄를 시도한 대 남성이 구속됐습니다서울 강남경찰서는 강간상...
4      서울연합뉴스대한약사회가 일부터 코로나바이러스 감염증 대응 체계를 사회적 거리두기에서...

...
4950     신종 코로나바이러스 감염증코로나사태 이후 가정의 달월에도 언택트비대면신품속도가 이어...
4951     는 소비자로부터 월 이용료 만만원을 받고 초고속 인터넷을 제공한다그런 브로드밴드가와...
4952     머리를 굽고 있는 오랑우탄몸을 굽는 행동을 따라 하는 것은 부정적 감정과 관련이 있...

```



```
4953     가 오는 일 정식 출시하는 스마트폰 벨벳이 사실상 공짜폰이 될 전망이다단말기 가격 ...
4954     이미지제공게티이미지뱅크 이미지제공게티이미지뱅크전자신문 전자신문인터넷무단전재 및 재배포...
Name: news, Length: 3987, dtype: object
```

```
1 len(df)
```

```
5249
```

```
1 print(df.isnull().sum())
```

```
news      0
code      0
dtype: int64
```

```
1 df.drop_duplicates(subset=['news'], inplace=True)
```

```
2 print('뉴스 기사의 갯수 :', len(df))
```

```
뉴스 기사의 갯수 : 3987
```

▼ 데이터 탐색

```
1 df['code'].value_counts().plot(kind='bar')
```

```

<matplotlib.axes._subplots.AxesSubplot at 0x7fe892c45cd0>/usr/local/lib/python3.7/dist-packages/matplotlib/backends/backend_agg.py:214:
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.7/dist-packages/matplotlib/backends/backend_agg.py:214:
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.7/dist-packages/matplotlib/backends/backend_agg.py:214:
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.7/dist-packages/matplotlib/backends/backend_agg.py:214:
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.7/dist-packages/matplotlib/backends/backend_agg.py:214:
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.7/dist-packages/matplotlib/backends/backend_agg.py:214:
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.7/dist-packages/matplotlib/backends/backend_agg.py:214:
font.set_text(s, 0.0, flags=flags)
/usr/local/lib/python3.7/dist-packages/matplotlib/backends/backend_agg.py:183:
font.set_text(s, 0, flags=flags)
/usr/local/lib/python3.7/dist-packages/matplotlib/backends/backend_agg.py:183:
font.set_text(s, 0, flags=flags)

```

```
1 print(df.groupby('code').size().reset_index(name='count'))
```

```

      code  count
0  IT/과학    903
1     사회   1670
2  생활/문화   1414

```

```

1 from konlpy.tag import Mecab
2 tokenizer = Mecab()

```

```

1 kor_text = '밤에 귀가하던 여성에게 범죄를 시도한 대 남성이 구속됐다서울 제주경찰서는 \
2             상해 혐의로 씨를 구속해 수사하고 있다고 일 밝혔다씨는 지난달 일 피해 여성을 \
3             인근 지하철 역에서부터 따라가 폭행을 시도하려다가 도망간 혐의를 받는다피해 \
4             여성이 저항하자 놀란 씨는 도망갔으며 신고를 받고 주변을 수색하던 경찰에 \
5             체포됐다피해 여성은 이 과정에서 경미한 부상을 입은 것으로 전해졌다'

```

```
800 | ██████████ ██████████ ██████████ |
```

```
1 print(tokenizer.morphs(kor_text))
```

```

['밤', '에', '귀가', '하', '던', '여성', '에게', '범죄', '를', '시도', '한', '대', '남성',
 | ██████████ ██████████ ██████████ |

```

▼ 불용어 제거

```
1 stopwords = ['에', '는', '은', '을', '했', '에게', '있', '이', '의', '하', '한', '다', '과', '때문', '할
```

```

1 def preprocessing(data):
2     text_data = []
3
4     for sentence in data:
5         temp_data = []
6         temp_data = tokenizer.morphs(sentence)
7         temp_data = [word for word in temp_data if not word in stopwords]

```

```

8         text_data.append(temp_data)
9
10    text_data = list(map(' '.join, text_data))
11    return text_data

```

```
1 text_data = preprocessing(df['news'])
```

```
1 print(text_data[0])
```

파주 시청 사진 제공 = 파주시 파주 시청 사진 제공 = 파주시 [파주 = 강근주 파 주 시 4 일 관내 취약

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.feature_extraction.text import CountVectorizer
3 from sklearn.feature_extraction.text import TfidfTransformer
4 from sklearn.naive_bayes import MultinomialNB
5 from sklearn import metrics
```

```
1 x_train, x_test, y_train, y_test = train_test_split(text_data, df['code'], random_state=42)
```

```
1 print('훈련용 뉴스 기사의 갯수 : ', len(x_train))
2 print('테스트용 뉴스 기사의 갯수 : ', len(x_test))
3 print('훈련용 레이블의 갯수 : ', len(y_train))
4 print('테스트용 레이블의 갯수 : ', len(y_test))
```

훈련용 뉴스 기사의 갯수 : 2990
 테스트용 뉴스 기사의 갯수 : 997
 훈련용 레이블의 갯수 : 2990
 테스트용 레이블의 갯수 : 997

```
1 # DTM
2 count_vect = CountVectorizer()
3 x_train_counts = count_vect.fit_transform(x_train)
```

```
1 # TF-IDF 행렬
2 tfidf_transformer = TfidfTransformer()
3 x_train_tfidf = tfidf_transformer.fit_transform(x_train_counts)
```

```
1 # 모델 학습 (나이브베이즈모델)
2 clf = MultinomialNB().fit(x_train_tfidf, y_train)
```

▼ 평가 및 테스트

```
1 def tfidf_vectorizer(data):
2     data_counts = count_vect.transform(data)
3     data_tfidf = tfidf_transformer.transform(data_counts)
4     return data_tfidf
```

[illegible]

```

1 new_sent = preprocessing(["인류는 결국에 저 세계·자구 심사 기준을 없애야 한다는 \
2                             주장이 나오는데 대해 “체계·자구 심사가 법안 지연의 수단으로 \
3                             쓰이는 것은 바람직하지 않다”면서도 “국회를 통과하는 법안 중 위헌\
4                             법률이 1년에 10건 넘게 나온다. 그런데 체계·자구 심사까지 없애면 매우
5 ]])

```

```

1 print(clf.predict(tfidf_vectorizer(new_sent)))

```

```
['사회']
```

```

1 new_sent = preprocessing(["인도 로맨틱 코미디 영화 <까립까립 싱글>(2017)을 봤을 때 나는 두 눈
2                             저 사람이 남자 주인공이라고? 노안에 가까운 이목구비와 기름때로 뭍힌 파
3                             대충 툭툭 던지는 말투 등 전혀 로맨틱하지 않은 외모였다. 반감이 일면서
4                             ‘난 외모지상주의자가 아니다’라고 자부했던 나에 대해 회의가 들었다.\
5                             티브이를 꺼버릴까? 다른 걸 볼까? 그런데, 이상하다. 왜 이렇게 매력 있
6                             같이 툭 불거진 눈망을 안에는 어떤 인도 배우에게서도 느끼지 못한 \
7                             부드러움과 선량함, 무엇보다 슬픔이 있었다. 2시간 뒤 영화가 끝나고 나

```

```

1 print(clf.predict(tfidf_vectorizer(new_sent)))

```

```
['생활/문화']
```

```

1 new_sent = preprocessing(["20분기 연속으로 적자에 시달리는 LG전자가 브랜드 이름부터 성능, 디자인
2                             적용한 LG 벨벳은 등장 전부터 온라인 커뮤니티를 뜨겁게 달궜다. 사용자들
3                             “슬림하다”는 반응을 보이며 LG 벨벳에 대한 기대감을 드러냈다.”])

```

```

1 print(clf.predict(tfidf_vectorizer(new_sent)))

```

```
['IT/과학']
```

```

1 y_pred = clf.predict(tfidf_vectorizer(x_test))
2 print(metrics.classification_report(y_test, y_pred))

```

	precision	recall	f1-score	support
IT/과학	0.93	0.79	0.85	240
사회	0.79	0.91	0.85	425
생활/문화	0.85	0.78	0.81	332
accuracy			0.84	997
macro avg	0.86	0.83	0.84	997
weighted avg	0.85	0.84	0.84	997

```
1
```

