

In [ ]:

```
#naver news crawling
from bs4 import BeautifulSoup
import urllib.request
```

In [ ]:

```
import numpy as np
import pandas as pd
```

In [ ]:

```
# URL 주소 [정치]
URL1 = np.array([])
for i in np.arange(100):
    num = 3053599 - i
    new_URL = np.array(['https://news.naver.com/main/read.nhn?
mode=LSD&mid=shm&sid1=100&oid=081&aid=000{}'.format(num)])
    URL1 = np.concatenate((URL1,new_URL))

# URL 주소 [경제]
URL2 = np.array([])
for i in np.arange(100):
    num = 2962982 - i
    new_URL = np.array(['https://news.naver.com/main/read.nhn?
mode=LSD&mid=shm&sid1=101&oid=025&aid=000{}'.format(num)])
    URL2 = np.concatenate((URL2,new_URL))

# URL 주소 [사회]
URL3 = np.array([])
for i in np.arange(100):
    num = 7876 - i
    new_URL = np.array(['https://news.naver.com/main/read.nhn?
mode=LSD&mid=shm&sid1=102&oid=629&aid=000000{}'.format(num)])
    URL3 = np.concatenate((URL3,new_URL))

# URL 주소 [문화]
URL4 = np.array([])
for i in np.arange(100):
    num = 4542990 - i
    new_URL = np.array(['https://news.naver.com/main/read.nhn?
mode=LSD&mid=shm&sid1=103&oid=018&aid=000{}'.format(num)])
    URL4 = np.concatenate((URL4,new_URL))

URL = np.array([URL1, URL2, URL3, URL4])
# 12/25 01:30 AM 기준
```

In [ ]:

```
URL1[0]
```

Out[ ]:

```
'https://news.naver.com/main/read.nhn?mode=LSD&mid=shm&sid1=100&oid=081&aid=0003053599'
```

In [ ]:

```
# 크롤링 함수
def get_text(URL):
    source_code_from_URL = urllib.request.urlopen(URL)
    soup = BeautifulSoup(source_code_from_URL, 'lxml', from_encoding='utf-8')
    text = ''
    for item in soup.find_all('div', id='articleBodyContents'):
        text = text + str(item.find_all(text=True))
    return text
```

In [ ]:

```
# 데이터 클리닝 함수
import re

# 클린 함수 (영어 및 특수 기호 제거)
def clean(before):
    after = re.sub('[a-zA-Z]', '', before)
    after = re.sub('[\{\}\[\]\./?.,;:\|\)*~`!^_-+<>@\#\$%\&\&=\(\)\'\"]', '', after)
    return after
```

In [ ]:

```
# 데이터 프레임으로 만들기
import pandas as pd

genre = ['정치', '경제', '사회', '문화']
text = np.array([])
category = np.array([])
for i in np.arange(4):
    for j in np.arange(100):
        cat = np.array([genre[i]])
        category = np.concatenate((category, cat))
        news = np.array([clean(get_text(URL[i][j]))])
        text = np.concatenate((text, news))

data0 = pd.DataFrame(np.array([text, category]), index=['text', 'category'])
data0 = data0.T
```

In [ ]:

```
from google.colab import drive
drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force\_remount=True).

In [ ]:

```
%pwd
```

Out [ ]:

```
'/content/drive/My Drive'
```

In [ ]:

```
%cd /content/drive/My Drive
```

```
/content/drive/My Drive
```

In [ ]:

```
# csv 파일로 저장
data0.to_csv("data0.csv")
```

In [ ]:

```
# csv 파일 불러오기
data0 = pd.read_csv('data0.csv')
```

In [ ]:

```
data0 = data0.loc[:, ['text', 'category']]
data0
```

Out [ ]:

text category



[ '가까운', '가장', '갈등', '강조했다', '같이', '걸어가야', '것이', '계속', '계정에', '공유하는', '과거의', '구독하기', '국가들과', '국은', '귀국길에', '그러나', '글을', '기자', '긴밀히', '나가야', '나라든', '네이버에서', '대통령과', '대통령은', '동안의', '때때로', '떠나며', '라는', '마치고', '모두', '무단전재', '무료만화', '문재인', '문화를', '미래지향적인', '발전시켜', '발전해', '밝혔다', '바침중', '보기', '분명히', '불거지는', '불행한', '사는', '샹그릴라호텔에서', '서울신문', '세기성', '세상화', '소셜네트웍서비스', '속에서', '수천', '쓰

```
참성', '아베', '악수하고', '앞서', '어느', '어울려', '언급했다', '없다', '역사로', '역사를', '역사와', '연합  
뉴스', '오랜', '오른', '올려', '요소가', '우리가', '우리는', '이날', '이라고', '이라며', '이라면서', '이런',  
'이렇게', '이웃', '인해', '일본', '일이', '일정을', '일현지시간', '있다', '자신의', '재배포금지', '정상회담에',  
'제목의', '중국', '직시하면서도', '청두', '청두를', '총리가', '최선을', '클릭', '하고', '한다', '한중일',  
'함께', '협력', '협력을', '협력해야', '출로']  
(1, 101)
```

In [ ]:

```
# 장르 별 최빈단어 10개씩 찾기
```

```
vocab1 = np.array([])  
for i in np.arange(100):  
    text = np.array([X[i]])  
    new_vocab = vectorizer.fit_transform(text)  
    vocab1 = np.concatenate((vocab1, vectorizer.get_feature_names()))
```

Out[ ]:

```
18176
```

In [ ]:

```
# train, test set 나누기
```

```
from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, stratify=y)
```

In [ ]:

```
# 단어를 분할하는 tokenizer - 안써도 될듯  
def tokenizer(text):  
    return text.split()
```

In [ ]:

```
#LogisticRegression
```

```
from sklearn.model_selection import GridSearchCV  
from sklearn.pipeline import Pipeline  
from sklearn.linear_model import LogisticRegression  
from sklearn.feature_extraction.text import TfidfVectorizer  
  
tfidf = TfidfVectorizer(strip_accents=None, lowercase=False, preprocessor=None)  
param_grid = [{'vect_ngram_range': [(1,1)], 'vect_tokenizer': [tokenizer, None], 'clf_penalty': ['l1',  
'l2'], 'clf_C': [1.0, 10.0, 100.0]},  
               {'vect_ngram_range': [(1,1)], 'vect_tokenizer': [tokenizer, None], 'vect_use_idf': [False],  
               'vect_norm': [None], 'clf_penalty': ['l1', 'l2'], 'clf_C': [1.0, 10.0, 100.0]}]  
lr_tfidf = Pipeline([('vect', tfidf), ('clf', LogisticRegression())])  
gs_lr_tfidf = GridSearchCV(lr_tfidf, param_grid, scoring='accuracy', cv=5, verbose=1, n_jobs=1)  
gs_lr_tfidf.fit(X_train, y_train)
```

In [ ]:

```
print(gs_lr_tfidf.best_params_)  
  
y_train_pred1 = gs_lr_tfidf.predict(X_train)  
y_test_pred1 = gs_lr_tfidf.predict(X_test)  
  
from sklearn.metrics import accuracy_score  
print(accuracy_score(y_train, y_train_pred1)) # train data에 대한 accuracy  
print(accuracy_score(y_test, y_test_pred1)) # test data에 대한 accuracy
```

```
{'clf_C': 1.0, 'clf_penalty': 'l1', 'vect_ngram_range': (1, 1), 'vect_tokenizer': <function to  
kenizer at 0x7faa73d44620>}  
1.0  
1.0
```

In [ ]:

```
# 분류 결과 logistic
from sklearn.metrics import confusion_matrix
print(confusion_matrix(y_test, y_test_pred1))
```

```
[[30  0  0  0]
 [ 0 30  0  0]
 [ 0  0 30  0]
 [ 0  0  0 30]]
```

In [ ]:

```
#KNN

from sklearn.model_selection import GridSearchCV
from sklearn.pipeline import Pipeline
from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.neighbors import KNeighborsClassifier
tfidf = TfidfVectorizer(strip_accents=None, lowercase=False, preprocessor=None)
param_grid = [{'vect__ngram_range': [(1,1)], 'vect__tokenizer': [tokenizer, None], 'knn__n_neighbors':
[3,5,10], 'knn__p': [1,2]},
               {'vect__ngram_range': [(1,1)], 'vect__tokenizer': [tokenizer, None], 'vect__use_idf': [
False],
               'vect__norm': [None], 'knn__n_neighbors': [3,5,10], 'knn__p': [1,2]}]
knn_tfidf = Pipeline([('vect', tfidf), ('knn', KNeighborsClassifier())])
gs_knn_tfidf = GridSearchCV(knn_tfidf, param_grid, scoring='accuracy', cv=5, verbose=1, n_jobs=1)
gs_knn_tfidf.fit(X_train, y_train)
```

Fitting 5 folds for each of 24 candidates, totalling 120 fits

[Parallel(n\_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.  
[Parallel(n\_jobs=1)]: Done 120 out of 120 | elapsed: 14.4s finished

Out[ ]:

```
GridSearchCV(cv=5, error_score='raise-deprecating',
             estimator=Pipeline(memory=None,
                                steps=[('vect',
                                         TfidfVectorizer(analyzer='word',
                                                           binary=False,
                                                           decode_error='strict',
                                                           dtype=<class 'numpy.float64'>,
                                                           encoding='utf-8',
                                                           input='content',
                                                           lowercase=False,
                                                           max_df=1.0,
                                                           max_features=None,
                                                           min_df=1,
                                                           ngram_range=(1, 1),
                                                           norm='l2',
                                                           preprocessor=None,
                                                           smooth_idf=True,
                                                           stop_word...
                                                           'vect__ngram_range': [(1, 1)],
                                                           'vect__tokenizer': [<function tokenizer at 0x7faa73d44620>,
                                                                               None]),
                                ('knn__n_neighbors': [3, 5, 10], 'knn__p': [1, 2],
                                'vect__ngram_range': [(1, 1)], 'vect__norm': [None],
                                'vect__tokenizer': [<function tokenizer at 0x7faa73d44620>,
                                                                               None],
                                'vect__use_idf': [False])],
             pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
             scoring='accuracy', verbose=1)
```

In [ ]:

```
print(gs_knn_tfidf.best_params_)
```

```
y_train_pred2 = gs_knn_tfidf.predict(X_train)
```

```
from sklearn.metrics import accuracy_score
print(accuracy_score(y_train, y_train_pred2)) # train data에 대한 accuracy
print(accuracy_score(y_test, y_test_pred2)) # test data에 대한 accuracy
```

In [ ]:

```
[[17  1  6  6]
 [ 0 27  0  3]
 [ 1  1 26  2]
 [ 0  2  2 26]]
```

```
#SVM

from sklearn.model_selection import GridSearchCV
from sklearn.pipeline import Pipeline
from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.svm import SVC

tfidf = TfidfVectorizer(strip_accents=None, lowercase=False, preprocessor=None)

param_grid = [{'vect_ngram_range': [(1,1)], 'vect_tokenizer': [tokenizer, None], 'svc_C': [1e3, 5e3, 1e4, 5e4, 1e5],
               'svc_gamma': [0.005, 0.01, 0.1], 'svc_kernel': ['poly', 'rbf', 'sigmoid']],
              {'vect_ngram_range': [(1,1)], 'vect_tokenizer': [tokenizer, None], 'vect_use_idf': [False],
               'vect_norm': [None], 'svc_C': [1e3, 5e3, 1e4, 5e4, 1e5],
               'svc_gamma': [0.005, 0.01, 0.1], 'svc_kernel': ['poly', 'rbf', 'sigmoid']}]

svc_tfidf = Pipeline([('vect', tfidf), ('svc', SVC())])
gs_svc_tfidf = GridSearchCV(svc_tfidf, param_grid, scoring='accuracy', cv=5, verbose=1, n_jobs=1)
gs_svc_tfidf.fit(X_train, y_train)
```

```
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.  
[Parallel(n_jobs=1)]: Done 900 out of 900 | elapsed: 2.8min finished
```

```
GridSearchCV(cv=5, error_score='raise-deprecating',
             estimator=Pipeline(memory=None,
                                steps=[('vect',
                                         TfIdfVectorizer(analyzer='word',
                                                           binary=False,
                                                           decode_error='strict',
                                                           dtype=<class 'numpy.float64'>,
                                                           encoding='utf-8',
                                                           input='content',
                                                           lowercase=False,
                                                           max_df=1.0,
                                                           max_features=None,
                                                           min_df=1,
                                                           ngram_range=(1, 1),
                                                           norm='l2',
                                                           preprocessor=None,
                                                           smooth_idf=True,
                                                           stop_word...
                                                           None)]),
             {'svc__C': [1000.0, 5000.0, 10000.0, 50000.0,
                          100000.0],
              'svc__gamma': [0.005, 0.01, 0.1]}).
```

In [ ]:

In [ ]:

In [ ]:

Out[ ]:

[illegible]



```

        lowercase=False,
        max_df=1.0,
        max_features=None,
        min_df=1,
        ngram_range=(1, 1),
        norm='l2',
        preprocessor=None,
        smooth_idf=True,
        stop_word...
    'vect__ngram_range': [(1, 1)],
    'vect__tokenizer': [<function tokenizer at 0x7faa73d44620>,
                        None]],
    {'rf__max_depth': [3, 5, 10],
     'rf__n_estimators': [50, 100, 200, 500],
     'vect__ngram_range': [(1, 1)], 'vect__norm': [None],
     'vect__tokenizer': [<function tokenizer at 0x7faa73d44620>,
                        None],
     'vect__use_idf': [False]}],
    pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
    scoring='accuracy', verbose=1)

```

In [ ]:

```

print(gs_rf_tfidf.best_params_)

y_train_pred4 = gs_rf_tfidf.predict(X_train)
y_test_pred4 = gs_rf_tfidf.predict(X_test)

from sklearn.metrics import accuracy_score
print(accuracy_score(y_train, y_train_pred4)) # train data에 대한 accuracy
print(accuracy_score(y_test, y_test_pred4)) # test data에 대한 accuracy

{'rf__max_depth': 3, 'rf__n_estimators': 50, 'vect__ngram_range': (1, 1), 'vect__tokenizer': <func
tion tokenizer at 0x7faa73d44620>}
1.0
1.0

```

In [ ]:

```

# 분류 결과 RandomForest
from sklearn.metrics import confusion_matrix
print(confusion_matrix(y_test, y_test_pred4))

[[30  0  0  0]
 [ 0 30  0  0]
 [ 0  0 30  0]
 [ 0  0  0 30]]

```