



LightGBM을 활용한 판결문 주장/사실/판단/결론 분류 모델

목표

- ▶ AIHub에서 텍스트 데이터 수집을 진행한다.
- ▶ 데이터에 대한 탐색적 분석을 수행하고 활용 목적에 알맞는 전처리를 수행한다.
- ▶ TF-IDF를 활용해 텍스트 데이터를 임베딩한다.
- ▶ 임베딩된 데이터로 LightGBM모델을 학습시켜 판결문의 주장/사실/판단/결론에 대한 분류를 수행한다.



진행 순서

1. 데이터 수집

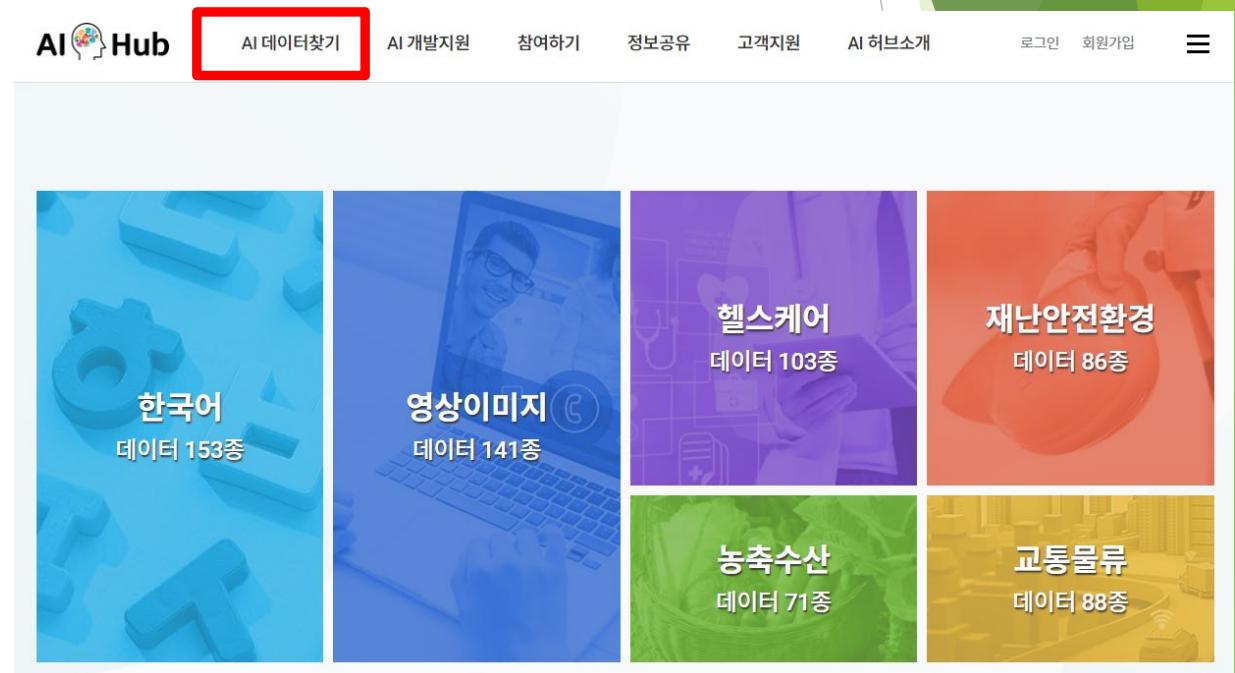
2. EDA 및 전처리

3. 모델링 및 결과 확인

1. 데이터 수집

1.1 AIHub에서의 데이터 수집

- [AIHub 사이트](#)에 접속한 뒤 상단의 AI 데이터찾기를 클릭한다.



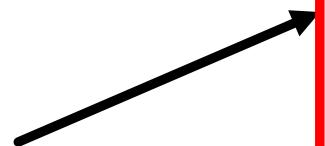
1. 데이터 수집

1.1 AIHub에서의 데이터 수집

1. 검색창에 사용하고자 하는 데이터를 검색하거나



2. 아래의 데이터 중 사용하고자 하는 데이터를 클릭한다.



The screenshot shows the AIHub 'Data Catalog' page. At the top, there is a search bar with the placeholder '검색어를 입력해주세요' and a blue button labeled '데이터셋 검색'. Below the search bar are two rows of filters: '분야 선택' (Field Selection) and '데이터유형 선택' (Data Type Selection). The '데이터유형 선택' row contains icons for Image, Video, Text, Audio, 3D, and Sensor. A red box highlights the search bar and the entire filter section. Below these filters, a list of datasets is displayed with a red box highlighting the last item. Each dataset entry includes a thumbnail, the dataset name, a 'BETA' badge, download counts, and a date range.

데이터셋	설명	버전	다운로드 수	등록일
전술 판정 영상 데이터(핸드볼)	BETA	2023-07	109	2022
전술 판정 영상 데이터(농구)	BETA	2023-07	111	2022
감정이 태깅된 자유대화 (청소년)	BETA	2023-07	0	2022

1. 데이터 수집

1.1 AIHub에서의 데이터 수집

- 본 실습에서는 법률/규정 (판결서, 약관 등) 텍스트 분석 데이터를 사용한다.

The screenshot shows the AIHub dataset page for '법률/규정 (판결서, 약관 등) 텍스트 분석 데이터'. The page includes a sidebar with a green decorative background, a header with a 'DATA' icon, and various filters like '#인공지능', '#리걸테크', '#판결문', '#학습용 데이터', and '#판례분석'. Below the header, there are tabs for '분야' (Law), '한국어' (Korean), '유형' (Type), and '텍스트' (Text). The main content area displays the dataset details: '갱신년월: 2023-05', '구축년도: 2021', '조회수: 7,298', '다운로드: 737', and '용량: 156.67 MB'. A prominent red box highlights the '다운로드' (Download) button. Other buttons include '샘플 데이터' (Sample Data) and a question mark icon. To the right, there's a '관심데이터 등록' (Register Interest) button with a count of '46'. At the bottom, a note says '※ 내국인만 데이터 신청이 가능합니다.' (Only domestic users can request data). The page also features sections for '데이터 개요' (Data Overview) and '데이터 변경이력' (Data Change History), which lists two versions: 1.1 (May 10, 2023) and 1.0 (July 29, 2022).

버전	일자	변경내용	비고
1.1	2023-05-10	원천데이터, 라벨링데이터, 샘플데이터 수정	
1.0	2022-07-29	데이터 최초 개방	

2. EDA 및 전처리

2.1 라이브러리 import 및 초기설정

- 사용하고자 하는 분석에 필요한 라이브러리들을 불러온다.

```
import pandas as pd
import numpy as np
import re
import io
import os
from tqdm.notebook import tqdm
tqdm.pandas()
from torch.utils.data import Dataset, DataLoader
from ml_things import fix_text
from sklearn.metrics import classification_report, accuracy_score

data_path = './archive/라벨링데이터/판결문'

# 분류 라벨
labels_ids = {'01_주장':0, '02_사실':1, '03_판단':2, '04_결론':3}

# How many labels are we using in training.
# This is used to decide size of classification head.
n_labels = len(labels_ids)
```

2. EDA 및 전처리

2.2 EDA 수행 및 데이터 전처리

- 사용하고자 하는 데이터 샘플을 불러와 살펴본다.

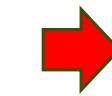
- json 확장자 파일에 9개의 key (info, concerned, org, disposal, mentionedItems, assrs, facts, dcss, close) 와 각 key에 해당하는 내용들이 포함되어 있다.

2. EDA 및 전처리

2.2 EDA 수행 및 데이터 전처리

1. 읽어온 JSON 객체에 특정 key의 존재 여부를 확인하기위한 `is_json_key_present` 함수를 정의한다.
2. 설정한 경로에서 데이터를 불러와 딕셔너리로 반환하는 `JsonDataset` 클래스를 정의한다.

```
# JSON 객체에서 특정 키(key)가 존재하는지 여부를 확인
def is_json_key_present(json, key):
    try:
        buf = json[key]
    except KeyError:
        return False
    return True
```



```
class JsonDataset(Dataset):

    def __init__(self, path):
        # Check if path exists.
        if not os.path.isdir(path):
            # Raise error if path is invalid.
            raise ValueError('Invalid `path` variable! Needs to be a directory')
        self.texts = []
        self.labels = []

        for label in ['01_민사', '02_형사', '03_행정']:
            sentiment_path = os.path.join(path, label)

            # Get all files from path.
            for root, directories, files in os.walk(sentiment_path):
                # Go through each file and read its content.
                for file_name in tqdm(files, desc=f'{root} files'):
                    file_path = os.path.join(root, file_name)
                    with open(file_path, "r") as json_file:
                        json_data = json.load(json_file)
                        if(is_json_key_present(json_data,"assrs") and is_json_key_present(json_data['assrs'], "acusrAssrs")) :
                            for data in json_data['assrs'][['acusrAssrs']] :
                                self.texts.append(fix_text(data))
                                self.labels.append('01_주장')

                        if(is_json_key_present(json_data,"assrs") and is_json_key_present(json_data['assrs'], "dedatAssrs")) :
                            for data in json_data['assrs'][['dedatAssrs']] :
                                self.texts.append(fix_text(data))
                                self.labels.append('01_주장')

                        if(is_json_key_present(json_data,"facts") and is_json_key_present(json_data['facts'], "bsisFacts")) :
                            for data in json_data['facts'][['bsisFacts']] :
                                self.texts.append(fix_text(data))
                                self.labels.append('02_사실')

                        if(is_json_key_present(json_data,"dcss") and is_json_key_present(json_data['dcss'], "countDcss")) :
                            for data in json_data['dcss'][['countDcss']] :
                                self.texts.append(fix_text(data))
                                self.labels.append('03_판단')

                        if(is_json_key_present(json_data,"close") and is_json_key_present(json_data['close'], "cnclsns")) :
                            for data in json_data['close'][['cnclsns']] :
                                self.texts.append(fix_text(data))
                                self.labels.append('04_결론')

        # Number of examples.
        self.n_examples = len(self.labels)
        return

    def __len__(self):
        return self.n_examples

    def __getitem__(self, item):
        return {'text':self.texts[item],
                'label':self.labels[item]}
```

2. EDA 및 전처리

2.2 EDA 수행 및 데이터 전처리

1. 사전 정의한 JSONdataset 클래스로 데이터를 불러온다.
2. 불러온 데이터에 `pd.DataFrame.from_records`를 적용해 데이터프레임으로 만든다.
3. `train_test_split()`을 통해 데이터프레임을 9:1의 비율로 나눈다.
4. label 별로 나누어져 있는 데이터프레임을 합쳐 훈련 데이터와 테스트 데이터를 만든다.

※ GridSearchCV의 cv argument를 통해 validation이 가능하므로 train, test로만 분리한다.

```
import pandas as pd
from sklearn.model_selection import train_test_split

# JSON Data
dataset = JsonDataset(path=data_path)
df = pd.DataFrame.from_records(dataset)
df.groupby('label').size()

df = df.sample(frac=0.05, random_state=123)

df1 = df[df['label'].str.startswith('01')]
df2 = df[df['label'].str.startswith('02')]
df3 = df[df['label'].str.startswith('03')]
df4 = df[df['label'].str.startswith('04')]

# 전체 데이터에서 train, test 분리 (9 : 1)
df1_train, df1_test = train_test_split(df1, test_size=0.1, random_state=123)
df2_train, df2_test = train_test_split(df2, test_size=0.1, random_state=123)
df3_train, df3_test = train_test_split(df3, test_size=0.1, random_state=123)
df4_train, df4_test = train_test_split(df4, test_size=0.1, random_state=123)

df_train = pd.concat([df1_train, df2_train, df3_train, df4_train])
df_test = pd.concat([df1_test, df2_test, df3_test, df4_test])
```

