

데이터 전처리 인공지능 데이터 전처리 결과서

□ 개요

- 산출물 단계 : 데이터 전처리
- 평가 산출물 : 인공지능 데이터 전처리 결과서
- 제출 일자 : 2025.02.27
- 깃허브 경로 : <https://github.com/SKNETWORKS-FAMILY-AICAMP/SKN06-FINAL-2Team>
- 작성 팀원 : 이세화

데이터 전처리 개요	<p>데이터 수집일자</p> <ul style="list-style-type: none">- 2025.02.12 (기존 크롤링 데이터를 벡터화 진행)- 2025.02.27 (추가 크롤링 진행 중, 이 데이터를 벡터화하여 최종적으로 사용 예정) <p>데이터 양</p> <p>[학습 data]</p> <ul style="list-style-type: none">- 웹툰: 23,585 개- 웹소설: 156,973 개- 테스트를 마무리하고 최종 테스트 단계부터 사용될 예정. 배포 시에는 이 데이터들 사용. <p>[테스트 data]</p> <ul style="list-style-type: none">- 프롬프트 엔지니어링을 효율적으로 진행하기 위해 sample 데이터를 사용 플랫폼(네이버, 네이버시리즈, 카카오톡, 카카오페이지)별로 각 장르에서 20개씩 가져온 데이터로 테스트 작업 진행.- 장르와 플랫폼들의 다양성을 잘 인식하는지 실험해보고자 데이터 양을 기준으로 나누기보단 각각에서 일부를 가져옴.- 벡터화 하여 chroma db에 저장.
---------------	---

전처리 도구

```
# JSON 데이터를 벡터 스토어에 저장하는 함수
def json_to_vector_store(json_filename, store_name):
    with open(json_filename, "r", encoding="utf-8") as f:
        total = json.load(f)

    documents = []
    for idx in range(len(total)):
        content = total[idx]
        merged_text = f"{content['title']}({content['platform']})-{content['description']}, 장르:{content['genre']}, 키워드: {content['keywords']}"
        if isinstance(content['keywords'], str):
            keywords = content['keywords']
        else:
            keywords = ", ".join(content['keywords'])
        documents.append(
            Document(
                page_content=merged_text,
                metadata={
                    "title": content["title"],
                    "platform": content["platform"],
                    "update_days": content["update_days"],
                    "genre": content["genre"],
                    "price": content["price_price"],
                    "episode": content["episode"],
                },
            )
        )
    return documents
```

데이터 추출 방식

[모델]

- “text-embedding-ada-002” 모델로 벡터화 진행.
- “text-embedding-3-small”, “text-embedding-ada-002”, “Bespinn SROBERTa” 3개의 모델로 각각 진행함.
그 중 “text-embedding-ada-002” 모델로 벡터화 했을 때 사용자 쿼리를 가장 잘 인식 해 이 모델로 진행.
- 1) “text-embedding-3-small”
query: “여주가 짝사랑하는 로맨스 웹툰 추천해줘.”
-> 주체(‘여주’), 장르(‘로맨스’) 인식, 종류(‘웹툰’) 을 잘 인식하지 못하고 짝사랑하는 내용이 담긴 일상물 혹은 웹소설을 추천함.
- 2) “Bespinn SROBERTa”
query: “여주가 짝사랑하는 로맨스 웹툰 추천해줘.”
-> 장르(‘로맨스’)를 가장 잘 인식하게 함, 그러나 로맨스 웹툰 중 아무 웹툰을 추천함.

불필요한 데이터 제거 기준

- 크게 2가지 기준에 해당하면 metadata에 넣지 않음

[1]

- llm 모델에 사용될 정보는 아니고, MySQL에서 정보로 담고 있어야 할 키 값들
 - 1) id
 - 2) type -> (각각의 벡터스토어가 존재해 사용 안 됨)
 - 3) thumbnail
 - 4) age_rating -> (회원정보의 나이로 나이에제한에 맞는 작품 추천으로 사용)
 - 5) url

[2]

- 사용자가 자주 질문하지 않아 넣지 않아도 된다고 판단한 부분(네이버 지식인 글을 토대로 검사함. 3.모델링 및 평가_수집된 데이터 및 전처리 문서_0차 6번을 참조)

- 1) author
- 2) illustrator
- 3) status

[3]

- 따로 자체 웹툰 평가 지표 score를 만들어 추출할 예정. score의 feature로 사용되는 데이터라 metadata에는 필요 없음.

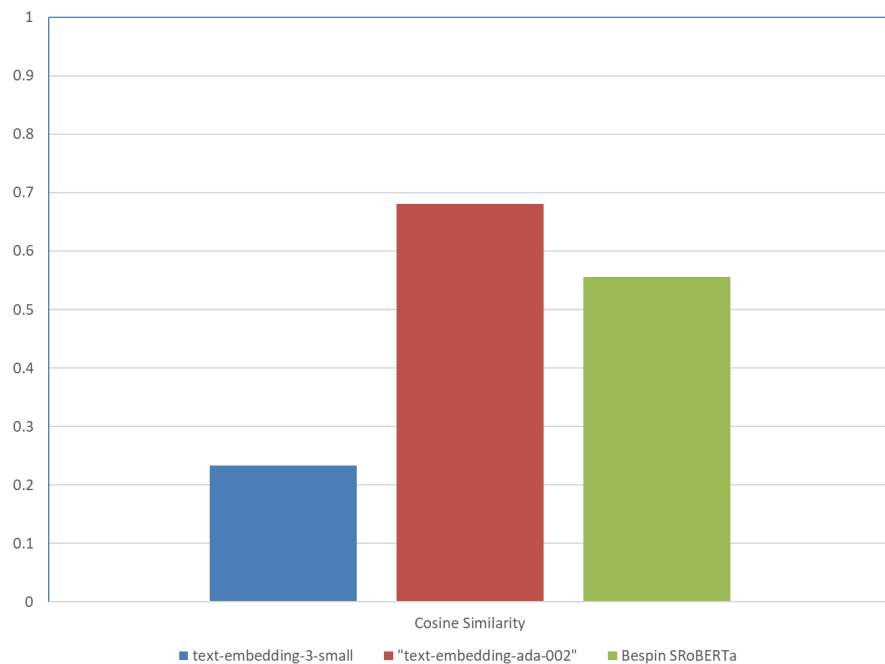
- 1) views
- 2) rating
- 3) likes
- 4) first_episode
- 5) comments_count
- 6) recent_comments_count

정제 방법

[한글-영어 벡터 유사도]

- 성능이 안 나오는 이유 중 하나가 사용자 쿼리에 쓰이는 단어들을 벡터 데이터베이스에서 잘 못 찾아서라고 가정함.
- 사용하고자 했던 모델의 한글-영어 벡터 유사도를 확인하고자 함. 자주 쓰이는 '장르'라는 단어와 'genre'의 벡터간 유사도(Cosine Similarity)를 살펴봄.

한국어/영어 벡터 유사도



- 유사도가 가장 높은 모델도 0.7 미만의 값을 보임.
- page_context 에 영어로 작성했던 부분을 모두 한국어로 바꿈.

- 1) genre -> '장르', keywords -> '키워드'
- 2) 이 정제 방법을 쓰고 훨씬 나은 답변을 내보임. 특히 장르를 잘 구분함.
- 3) 그러나 이렇게 바뀜으로서 이번 프로젝트에서는 한국어 사용자로 제한됨. 더 많은 시간을 투자해 다양한 언어 모델로 각각 사용이 가능하게끔 발전시키고 싶음.

```
documents = []
for idx in range(len(total)):
    content = total[idx]
    merged_text = f"title:{content['title']},{(content['platform'])},{content['description']},genre:{content['genre']}, keywords: {content['keywords']}"
    if isinstance(content['keywords'], str):
        keywords = content['keywords']
    else:
        keywords = content['keywords']

documents = []
for idx in range(len(total)):
    content = total[idx]
    merged_text = f"{content['title']}({content['platform']})-{content['description']}, 장르:{content['genre']}, 키워드: {content['keywords']}"
    if isinstance(content['keywords'], str):
        keywords = content['keywords']
    else:
        keywords = content['keywords']
```

[이름 변경]

- llm모델을 테스트하는 과정에서 줄거리를 더 잘 인식해야한다는 결론이 나와 기존 줄거리를 의미하던 'description' -> 'synopsis'로 변경
- '월 연재' -> '월요일'
 - 1) query: "수요일에 볼 만한 웹소설 추천해줘." -> 이런 경우 모든 모델에서 '수요일'을 전혀 인식하지 못하고 '월요일', '토요일'에 연재되는 작품을 추천함.
 - 2) 월 연재의 벡터값과 '월요일'의 벡터값 간의 유사도가 낮음이 원인으로 보고, 크롤링 데이터에서 연재 유무와 요일을 분리하고 '월요일'같은 형식으로 통일함.

[추후 예정]

- llm 모델 성능이 잘 나오냐? 안 나오냐? 성능에 따라 아래 전처리를 다양하게 시도 할 예정
 - 1) '?','!' 등의 특수 문자 제거
 - 2) 'synopsis'를 문장 단위로 / 다 합쳐진 page_context를 token수 단위로 chunking

데이터 전처리 결과

● 결고

[illegible]

- 향후 사용계획
 - 현재 위 데이터로 llm 모델 테스트 작업 진행 중.
 - 앞으로도 목표 성능이 나올 때까지 데이터 전처리&프롬프트 엔지니어링을 진행 예정.
 - 모델 작업이 완료되면 total_data 동일하게 백터화하여 최종 배포에 사용.