

# Assignment 6

## 프로젝트 최종 보고서



과목명	영어분석을 위한 기계학습
담당교수	최승택
제출일	2025년 12월 11일
전공	한국어교육과
학번	202302592
이름	임채운

## 1. 프로젝트 개요

### 1.1. 프로젝트명

: 교재 학습활동 추출 도구

### 1.2. 문제 상황

사범대 전공 특성상 교과서나 교재를 보아야 하는 경우가 매우 많다. 교과서의 종류가 과목별, 수준별, 학년별로 전부 다르고 출판사도 다양하여 교재 분석을 하거나 수업 시연 준비를 위해 교과서를 검토하는 데에 시간이 많이 소요된다. 이 수고를 덜기 위해 교과서 안에 어떤 글이 들어있는지, 어떤 질문이 들어있는지 등을 추출하여 보여주는 도구가 필요하다고 생각했다.

## 2. 진행 과정

### 2.1. 주제 설정 및 문제 정의

앞선 문제 상황을 해결하기 위해, 교과서 속 학습활동 질문 추출 도구를 개발하는 것을 이 프로젝트의 주제로 설정하였다. 처음에는 ‘성취 기준’이나 ‘작품’을 기준으로 학습활동을 추출하는 것을 목표로 하여 데이터 수집과 분석까지 진행하였으나, 실제 모델을 학습하는 과정이 너무 복잡하고 이 계획대로 프로젝트를 진행하기에는 개인 역량이 부족하다고 판단하여 기준을 설정하지 않고 단순히 지시문만 추출할 수 있는 도구를 개발하는 것으로 계획을 일부 수정하였다.

### 2.2. 데이터 수집 및 분석

#### - 최초 계획에 따른 데이터 수집 및 분석

데이터 수집과 분석 과정에서는 최초 계획에 따라 Assignment 4에서 ‘성취 기준’을 기준으로 설정하여 교과서 속 지시문을 수동으로 수집하였다. 수집한 데이터에 대한 개요는 다음과 같다.

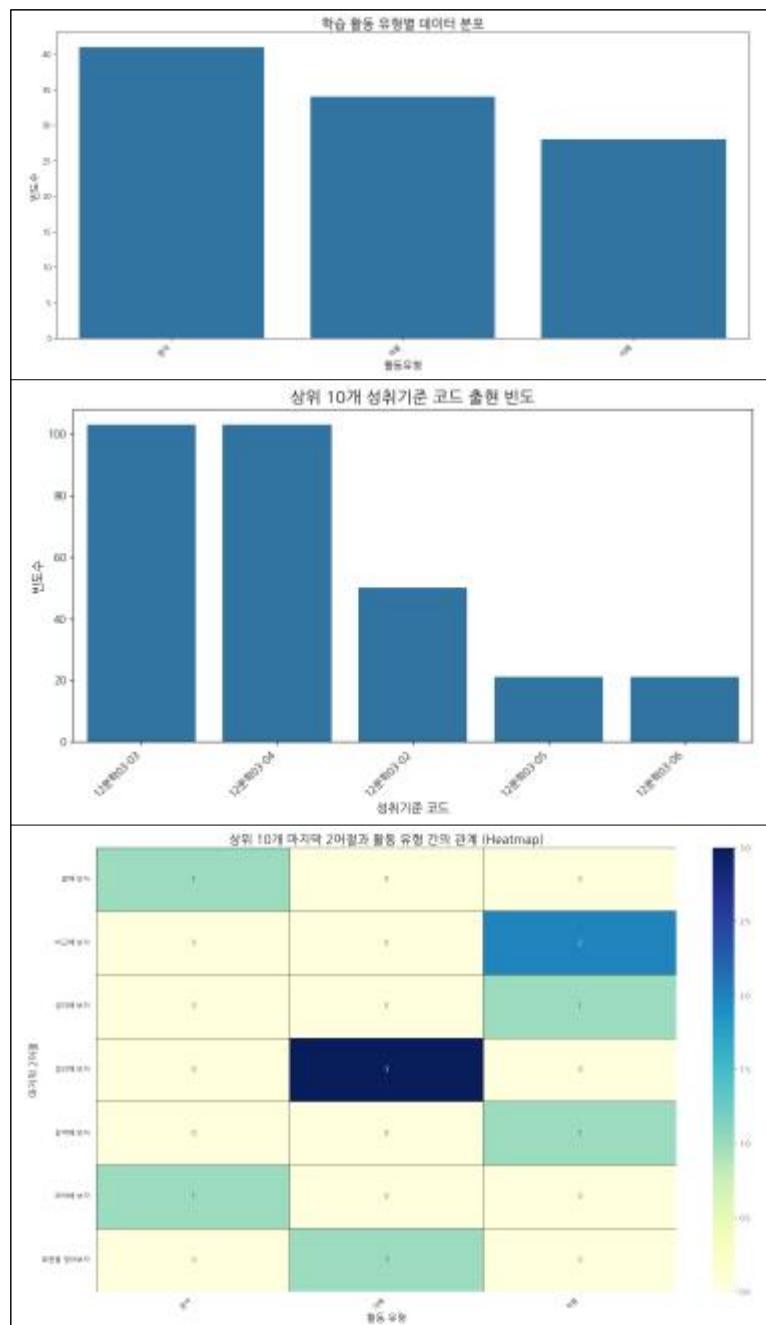
수집 데이터	2015 국어과 교육과정 성취기준, 국어 문학 교과서 (천재교과서, 동아출판, 비상교육, 천재교육)
개수	교육과정 성취기준 16개, 교과서 학습활동 지시문 103개

교육과정 성취기준은 ‘성취기준 코드 / 영역 / 학년 / 내용’으로 나누어 정리하였고, 교과서 학습활동 지시문은 ‘출판사 / 단원명 / 페이지 / 코드 / 글감제목 / 글쓴이 / 학습활동 지시문 / 활동 유형 / 질문 형식’으로 나누어 정리하였다. 수집한 자료들은 모두 csv 파일 형태로 저장하였다.

데이터 분석 과정에서 결측값은 존재하지 않았고, 활동 유형별 분포 확인 결과 ‘분석’, ‘적용’, ‘이해’ 순으로 많았던 것을 확인하였다. 이 유형은 학습활동 순서대로 이해, 분석, 적용 영역으로 임의로 부여한 것으로, 라벨링하는 주체에 따라 그 범위가 조금씩 달라질 수 있다. 성취 기준 코드의 출현 빈도 확인 결과 출현 빈도로는 ‘12문학03-03’과 ‘12문학03-04’가 가장 많았고, 이 두 항목은 수집된 모든 데이터에 적용된 코드임을 확인하였다. 문학 교과서에

한정되고, 특히 '문학사' 소단원에 한정하여 데이터를 수집하고 정리하였기 때문에 이 결과가 나타난 것으로 보인다. 이에 관해서는 향후 불균형한 데이터를 보완할 필요가 있다. 마지막으로는 학습활동 질문과 활동 유형과의 관계를 찾아보고자 하였다. 이에 따라, '학습활동 질문'의 마지막 2어절을 추출하여 활동 유형 '이해, 분석, 적용'과의 관계를 파악해 보았다. (예시: '말해 보자'로 끝나는 질문은 '이해' 유형에 해당할 것이다) 그러나 분석 결과 두 항목 간의 관계는 나타나지 않았고, 이 질문을 통해 활동 유형을 파악하기는 어려울 것으로 예상되었다.

다음은 데이터 분석과 관련된 그래프다.



수집한 데이터가 문학 교과서의 문학사 부분에만 해당하여 데이터의 불균형이 나타난 것으로 보이므로, 이를 보완하기 위해 국어의 다른 세부 과목 데이터나 다른 소단원의 데이터를 더 수집하고자 하였다.

#### - 수정된 계획에 따른 데이터 수집

계획이 수정됨에 따라 불필요한 성취 기준 데이터는 폐기하고, 기존에 수집한 학습활동 지시문만 채택하여 다시 정리하였다. 이 과정에서 데이터를 조금 더 추가하여 지시문 데이터를 총 115개를 수집하였다. 지시문이 아닌 일반 문장의 경우, 매우 정제된 패턴을 보이는 교과서의 문장에서 벗어나서 더 다양한 형태의 문장을 수집하기 위해 한국어 위키피디아에서 193개의 문장을 수집하였다. 지시문은 '1'을 할당하였고, 일반 문장은 '0'을 할당하여 정리한 자료는 csv 파일 형태로 저장하였다. 다음은 라벨링한 데이터의 예시다.

sentence	label
단군의 탄생 과정을 중심으로 '변화'라는 관점에서 살펴보고 빈칸을 채워 보자.	1
이 대회의 예선 방식은 조별 리그와 플레이오프로 되어 있다.	0

계획 수정에 따라 시간 관계상 이 데이터에 대해서는 분석 절차를 생략하였다.

### 2.3. ML 모델 학습 및 평가

한국어 데이터 처리를 위해 KoBERT 기반 정보추출 모델(skt/kobert-base-v1)을 활용하여 Assignment 4에서 수집한 데이터를 바탕으로 모델을 학습하였다. 다음은 학습 과정에서 분할한 데이터의 개수다.

분류	개수(개)
Train	196
Validation	50
Test	62

이를 바탕으로 학습을 시작하였고, 최초 epoch을 3으로 하여 실행하였으나, evaluation과 inference 성능이 너무 좋지 않아 epoch의 수를 10으로 늘렸다. epoch을 3에서 10으로 점차 늘린 결과, epoch의 수를 늘릴수록 validation accuracy가 항상 높아지는 것은 아님을 발견하였다. 그러나 전반적으로 우상향하는 추세를 보이고 있고, 더불어 training loss가 함께 줄어들고 있는 모습을 볼 수 있었다. epoch 10일 때의 주요 성능지표를 살펴보면, 지시문 추출 기준으로 Accuracy가 0.8, Precision 1.0, Recall 0.65, F1-score 0.79였다.

epoch을 늘릴수록 성능은 개선되었지만, Recall 수치가 너무 낮았다. 이 도구의 활용 목적이 '지시문 추출 후 검토'이므로 지시문이 아닌 문장을 추출하더라도 지시문인 문장을 빠짐없이 추출하도록 하고자 Recall 수치를 높이는 방식으로 조정하였다. 이에 따라, learning rate 를 2e-5에서 3e-5로 조정하여 모델 학습을 진행하였고, recall 수치가 0.77, F1-score가 0.87로 개선되었다. 다음 표는 모델 학습 및 평가에 관한 요약이다.

— 3단계: Test set 최종 평가 수행 — 평가 지표 (Precision, Recall, F1-score, Support) :				
	precision	recall	f1-score	support
0	0.8000	1.0000	0.8889	36
1	1.0000	0.6538	0.7907	26
accuracy			0.8548	62
macro avg	0.9000	0.8269	0.8396	62
weighted avg	0.8839	0.8548	0.8477	62

—————  
<learning rate 조정 전>

— 3단계: Test set 최종 평가 수행 — 평가 지표 (Precision, Recall, F1-score, Support) :				
	precision	recall	f1-score	support
0	0.8571	1.0000	0.9231	36
1	1.0000	0.7692	0.8696	26
accuracy			0.9032	62
macro avg	0.9296	0.8846	0.8963	62
weighted avg	0.9171	0.9032	0.9005	62

—————  
<learning rate 조정 후 최종>

### 3. 모델을 서비스로 만든 구조

#### 3.1. 서비스화

학습을 완료하고 실제 사용이 가능하도록 CLI 기반 도구로 만들었다. txb\_inference.py 파일이 핵심 로직이며, 모델 가중치는 instruction\_classifier\_model 폴더에 저장하였고, 모델 서비스에 필요한 라이브러리 목록을 requirements.txt에 저장하였다.

python 기반 서비스이기 때문에 이 서비스를 사용하기 위해서는 python이 설치된 환경이 필요하다. 또한, txt 파일을 입력으로 받기 때문에, txt 형태가 아닌 다른 형태의 파일이라면 미리 변환이 필요하다.

#### 3.2. 시스템 구조

- 1) CMD나 Anaconda prompt에서 사용자로부터 txt 파일의 경로를 입력받는다.

실행형식:

```
python      txb_inference.py      input_data\{파일명}.txt      -model-dir
instruction_classifier_model
```

- 2) 입력받은 파일에서 모델이 사전 학습된 ‘지시문’만 추출한다.
- 3) 추출된 지시문들이 콘솔 창을 통해 사용자에게 제시된다.

### 4. 실제 사용 결과

#### 4.1. 5회 이상 사용 기록 (스크린샷)

- 1) <국어논리및논술> 국어 교과서 제작에 필요한 교재 분석

<국어논리및논술> 과목에서 논리 교과서 제작을 하였다. 이를 위해 중학교 국어 교과서 속 논리 논술 단원을 분석하였고, 그 과정에서 이 도구를 활용하여 교과서 5종의 학습활동 속 지시문을 추출하여 교과서 제작에 참고하였다. OCR을 활용하여 미리 교과서에서 수동 추출한 텍스트를 txt 파일로 저장하여 도구에 제공하였고, 그 파일을 바탕으로 도구가 추출한 문장들을 출력하였다. 전반적으로 지시문을 거의 빠짐없이 추출하는 편이었으나, 불필요한 단어나 지시문이 아닌 문장들을 함께 추출하는 모습을 보였다. 다음은 실제 도구를 사용한 스크린샷 기록이다.

<pre>Anaconda Prompt - conda nv + ~ (robert_final_env) C:\final&gt;python txb_inference.py input_data\2015_miraen.txt --model_dir instruction_classifier_model 모델 로드 중: instruction_classifier_model [모델 로드 성공.]  [ 최종 주제 결과 ] 총 문장 수: 186 주술된 문서 수: 24  1. 2 자신이 노을이라면 어떤 방법으로 이름이를 설득할 수 있을지 생각해 보자.   2. (1) 이 글의 대목과 같이 행위할 때, 반란에 불미한 말맞은 말을 써 보자.   3. 그나마,   (3) 위 활동을 내용으로 하여 이 글의 구조를 파악해 보자.   4. 그나마   이제 최근 내 민족 비스 2 다음 글에 쓰인 논증 방법을 파악해 보자.   9. (1) 이 글의 문단은 이가 주장하는 바가 무엇인지 말해 보자.   10. (2) 다음 글을 참고하여 이 글의 논증 방법과 구조를 파악해 보자.   7. 논증 방법을 중심으로 글의 구조를 말 아보고, 논증의 대상성을 판단해 보자.   8. (1) 이 글의 글은 이가 주장하는 바가 무엇인지 말해 보자.   9. (2) 다음 글을 참고하여 이 글의 논증 방법과 구조를 파악해 보자.   10. 논증 방법을 중심으로 글의 구조를 말 아보고, 논증의 대상성을 판단해 보자.   11. 또한, 논증에서 근거를 '전제'라고도 하고, 주장은 '결론'이</pre>	<pre>Anaconda Prompt - conda nv + ~ 18. 논증 방법을 중심으로 글의 구조를 말 아보고, 논증의 대상성을 판단해 보자.   11. 또한, 논증에서 근거를 '전제'라고도 하고, 주장은 '결론'이라고도 한다는 점을 다시 한번 짚어보세요.   12. = 이끌기 논증의 특장성을 든다할 때.   (3) 이 글에서 주장은 떨치는 방식이 타당한지 판단해 보자.   13. (1) 논증방법 파악하며 유기 185   4 단원 문제해결 작성2/ 이끌기 논증방법을 사용하여 주장할 때는 전개 가운데 일부가 생략되기도 해요. 14. 설득의 힘 1 다음 글에 쓰인 논증 방법을 파악해 보자.   15. (3) 논증방법 파악하며 유기 185   4 단원 문제해결 작성2/ 2 정반대에서 모음을과 함께 논증 방법을 사용하여 이야기를 나누어 보자.   16. 예가 활용으로   좋은 일과를 보자.   17. 고정형 우리가 생활하는 중요한 공간에 다.   18. 설득의 힘</pre>
<pre>Anaconda Prompt - conda nv + ~  좋은 영화를 보자.   17. 교실은 우리가 생활하는 중요한 공간이다.   18. 설득의 힘 매운 내용을 떠올리며 반찬에 들어갈 알맞은 말을 찾아 보자.   19. (1)   은/는 일반적 풀리나 법칙에서 구체적 사실을 만들어 내는 논증 방법으로, 짧은 대전제를 사용하는 남다른 병법이다.   20. (3) 등 여성의 대상이 비슷한 속성을 가진다는 것을 근거로 다른 속성도 유사성을 것으로 추론하는 것을 (이)라고 한다.   21. 2 다음을 위고 몇몇은 말은 이에, 몇은 말은 3에 표시해 보자.   22. (1) 논증이란 근거를 들어 주장이 합리를 증명하는 것으로, 근거가 참마연 주장을 언제나 참이다.   23. (Q, X)   (3) 논증 방법을 파악하며 글을 읽으면 글의 타당성을 판단하는데 도움이 된다.   24. 내가 확실히 알고 있던 것에 나만의 편안 편견은 아닌지 의심해 본다.</pre>	
<pre>Anaconda Prompt - conda nv + ~ (robert_final_env) C:\final&gt;python txb_inference.py input_data\2022_miraen.txt --model_dir instruction_classifier_model 모델 로드 중: instruction_classifier_model [모델 로드 성공.]  [ 최종 주제 결과 ] 총 문장 수: 228 주술된 문서 수: 31    1. 행동이가 발견한 문제 상황을 써 보자.   2. 개념 쪽   3. 행동이와 친구들이 글의 주장과 이유, 예상 목표를 정리해 보자.   4. 그나마 우리가 글을 읽을 때 주제와 주장을 보면 찾는 이유를 정리해 보자.   5. 행동이와 친구들이 수집한 자료를 살펴보며, 자료를 수집할 때 유의할 점을 생각해 보자.   6. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   7. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   8. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   9. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   10. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   11. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   12. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   13. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   14. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   15. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   16. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   17. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   18. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   19. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   20. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   21. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   22. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   23. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   24. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   25. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   26. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   27. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   28. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.   29. 행동이와 친구들이 물질의 어떤 성분에 영향을 미친다고 추론하는 경우에 필요한 커뮤니케이션의 흐름을 알아보자.</pre>	<pre>Anaconda Prompt - conda nv + ~   5. 가짜연설 때 이   6. 가짜연설은 오직 물에   7. 어떤 연설에 있을지 조사해   8. 모리고 책을 읽었어.   9. 단상문트에 가짜연설 척기되   10. 면 단망과 선망을 조절해 주고   11. 더 부드러운 말을 만들어 준다.   12. 다른 풀과 기준에 따라 실률이 미와 친구들이 작성한 개요를 접검   13. 예제에 보자.   14. 자료 8   15. 행동이와 친구들이 정리한 자료 활용 방안을 바탕으로 하여 근거를 마련해 보자.   16. 예유 1. 카페인을 많이 섭취하면 신체에 부정적인 영향을 미   17. 글 정리 내용을 조작하고 주장의 분명하게   18. 드러나는 행동으로 네거티브로.   19. 3. 행동이와 친구들이 고친 것 외에 수절할 부분이 있다면 이   20. 예제에 보자.   21. 4. 다음은 근거를 들고 적절한 표현을 사용하여   22. 주장하는 글을 짜볼까?   23. 5. 우리 주변에서 청소년과 관련한 문제 상황을 짚어보자.   24. 6. 행위 속에서   25. 7. 주장은 되번 힘을 다해 한글 문서로 만들면 좋다.   26. 8. 2 1에서 떠올린 문제 상황 가운데 하나를 골라 자신의 주장   27. 9. 예상, 예상 목록을 정리해 보자.</pre>
<pre>Anaconda Prompt - conda nv + ~   10. (1)주장하는 글의 내용을 조작하기 위한 개요를 작성해 보자.   20. (2) 다음 풀과 기준에 따라 개요를 정리해 보자.   21. 개요   22. 개요   23. 개요   24. 개요   25. 개요   26. 개요   27. 개요   28. 개요   29. 개요   30. 개요   31. 개요   32. 개요   33. 개요   34. 개요   35. 개요   36. 개요   37. 개요   38. 개요   39. 개요   40. 개요   41. 개요   42. 개요   43. 개요   44. 개요   45. 개요   46. 개요   47. 개요   48. 개요   49. 개요   50. 개요   51. 개요   52. 개요   53. 개요   54. 개요   55. 개요   56. 개요   57. 개요   58. 개요   59. 개요   60. 개요   61. 개요   62. 개요   63. 개요   64. 개요   65. 개요   66. 개요   67. 개요   68. 개요   69. 개요   70. 개요   71. 개요   72. 개요   73. 개요   74. 개요   75. 개요   76. 개요   77. 개요   78. 개요   79. 개요   80. 개요   81. 개요   82. 개요   83. 개요   84. 개요   85. 개요   86. 개요   87. 개요   88. 개요   89. 개요   90. 개요   91. 개요   92. 개요   93. 개요   94. 개요   95. 개요   96. 개요   97. 개요   98. 개요   99. 개요   100. 개요   101. 개요   102. 개요   103. 개요   104. 개요   105. 개요   106. 개요   107. 개요   108. 개요   109. 개요   110. 개요   111. 개요   112. 개요   113. 개요   114. 개요   115. 개요   116. 개요   117. 개요   118. 개요   119. 개요   120. 개요   121. 개요   122. 개요   123. 개요   124. 개요   125. 개요   126. 개요   127. 개요   128. 개요   129. 개요   130. 개요   131. 개요   132. 개요   133. 개요   134. 개요   135. 개요   136. 개요   137. 개요   138. 개요   139. 개요   140. 개요   141. 개요   142. 개요   143. 개요   144. 개요   145. 개요   146. 개요   147. 개요   148. 개요   149. 개요   150. 개요   151. 개요   152. 개요   153. 개요   154. 개요   155. 개요   156. 개요   157. 개요   158. 개요   159. 개요   160. 개요   161. 개요   162. 개요   163. 개요   164. 개요   165. 개요   166. 개요   167. 개요   168. 개요   169. 개요   170. 개요   171. 개요   172. 개요   173. 개요   174. 개요   175. 개요   176. 개요   177. 개요   178. 개요   179. 개요   180. 개요   181. 개요   182. 개요   183. 개요   184. 개요   185. 개요   186. 개요   187. 개요   188. 개요   189. 개요   190. 개요   191. 개요   192. 개요   193. 개요   194. 개요   195. 개요   196. 개요   197. 개요   198. 개요   199. 개요   200. 개요   201. 개요   202. 개요   203. 개요   204. 개요   205. 개요   206. 개요   207. 개요   208. 개요   209. 개요   210. 개요   211. 개요   212. 개요   213. 개요   214. 개요   215. 개요   216. 개요   217. 개요   218. 개요   219. 개요   220. 개요   221. 개요   222. 개요   223. 개요   224. 개요   225. 개요   226. 개요   227. 개요   228. 개요   229. 개요   230. 개요   231. 개요   232. 개요   233. 개요   234. 개요   235. 개요   236. 개요   237. 개요   238. 개요   239. 개요   240. 개요   241. 개요   242. 개요   243. 개요   244. 개요   245. 개요   246. 개요   247. 개요   248. 개요   249. 개요   250. 개요   251. 개요   252. 개요   253. 개요   254. 개요   255. 개요   256. 개요   257. 개요   258. 개요   259. 개요   260. 개요   261. 개요   262. 개요   263. 개요   264. 개요   265. 개요   266. 개요   267. 개요   268. 개요   269. 개요   270. 개요   271. 개요   272. 개요   273. 개요   274. 개요   275. 개요   276. 개요   277. 개요   278. 개요   279. 개요   280. 개요   281. 개요   282. 개요   283. 개요   284. 개요   285. 개요   286. 개요   287. 개요   288. 개요   289. 개요   290. 개요   291. 개요   292. 개요   293. 개요   294. 개요   295. 개요   296. 개요   297. 개요   298. 개요   299. 개요   300. 개요   301. 개요   302. 개요   303. 개요   304. 개요   305. 개요   306. 개요   307. 개요   308. 개요   309. 개요   310. 개요   311. 개요   312. 개요   313. 개요   314. 개요   315. 개요   316. 개요   317. 개요   318. 개요   319. 개요   320. 개요   321. 개요   322. 개요   323. 개요   324. 개요   325. 개요   326. 개요   327. 개요   328. 개요   329. 개요   330. 개요   331. 개요   332. 개요   333. 개요   334. 개요   335. 개요   336. 개요   337. 개요   338. 개요   339. 개요   340. 개요   341. 개요   342. 개요   343. 개요   344. 개요   345. 개요   346. 개요   347. 개요   348. 개요   349. 개요   350. 개요   351. 개요   352. 개요   353. 개요   354. 개요   355. 개요   356. 개요   357. 개요   358. 개요   359. 개요   360. 개요   361. 개요   362. 개요   363. 개요   364. 개요   365. 개요   366. 개요   367. 개요   368. 개요   369. 개요   370. 개요   371. 개요   372. 개요   373. 개요   374. 개요   375. 개요   376. 개요   377. 개요   378. 개요   379. 개요   380. 개요   381. 개요   382. 개요   383. 개요   384. 개요   385. 개요   386. 개요   387. 개요   388. 개요   389. 개요   390. 개요   391. 개요   392. 개요   393. 개요   394. 개요   395. 개요   396. 개요   397. 개요   398. 개요   399. 개요   400. 개요   401. 개요   402. 개요   403. 개요   404. 개요   405. 개요   406. 개요   407. 개요   408. 개요   409. 개요   410. 개요   411. 개요   412. 개요   413. 개요   414. 개요   415. 개요   416. 개요   417. 개요   418. 개요   419. 개요   420. 개요   421. 개요   422. 개요   423. 개요   424. 개요   425. 개요   426. 개요   427. 개요   428. 개요   429. 개요   430. 개요   431. 개요   432. 개요   433. 개요   434. 개요   435. 개요   436. 개요   437. 개요   438. 개요   439. 개요   440. 개요   441. 개요   442. 개요   443. 개요   444. 개요   445. 개요   446. 개요   447. 개요   448. 개요   449. 개요   450. 개요   451. 개요   452. 개요   453. 개요   454. 개요   455. 개요   456. 개요   457. 개요   458. 개요   459. 개요   460. 개요   461. 개요   462. 개요   463. 개요   464. 개요   465. 개요   466. 개요   467. 개요   468. 개요   469. 개요   470. 개요   471. 개요   472. 개요   473. 개요   474. 개요   475. 개요   476. 개요   477. 개요   478. 개요   479. 개요   480. 개요   481. 개요   482. 개요   483. 개요   484. 개요   485. 개요   486. 개요   487. 개요   488. 개요   489. 개요   490. 개요   491. 개요   492. 개요   493. 개요   494. 개요   495. 개요   496. 개요   497. 개요   498. 개요   499. 개요   500. 개요   501. 개요   502. 개요   503. 개요   504. 개요   505. 개요   506. 개요   507. 개요   508. 개요   509. 개요   510. 개요   511. 개요   512. 개요   513. 개요   514. 개요   515. 개요   516. 개요   517. 개요   518. 개요   519. 개요   520. 개요   521. 개요   522. 개요   523. 개요   524. 개요   525. 개요   526. 개요   527. 개요   528. 개요   529. 개요   530. 개요   531. 개요   532. 개요   533. 개요   534. 개요   535. 개요   536. 개요   537. 개요   538. 개요   539. 개요   540. 개요   541. 개요   542. 개요   543. 개요   544. 개요   545. 개요   546. 개요   547. 개요   548. 개요   549. 개요   550. 개요   551. 개요   552. 개요   553. 개요   554. 개요   555. 개요   556. 개요   557. 개요   558. 개요   559. 개요   560. 개요   561. 개요   562. 개요   563. 개요   564. 개요   565. 개요   566. 개요   567. 개요   568. 개요   569. 개요   570. 개요   571. 개요   572. 개요   573. 개요   574. 개요   575. 개요   576. 개요   577. 개요   578. 개요   579. 개요   580. 개요   581. 개요   582. 개요   583. 개요   584. 개요   585. 개요   586. 개요   587. 개요   588. 개요   589. 개요   590. 개요   591. 개요   592. 개요   593. 개요   594. 개요   595. 개요   596. 개요   597. 개요   598. 개요   599. 개요   600. 개요   601. 개요   602. 개요   603. 개요   604. 개요   605. 개요   606. 개요   607. 개요   608. 개요   609. 개요   610. 개요   611. 개요   612. 개요   613. 개요   614. 개요   615. 개요   616. 개요   617. 개요   618. 개요   619. 개요   620. 개요   621. 개요   622. 개요   623. 개요   624. 개요   625. 개요   626. 개요   627. 개요   628. 개요   629. 개요   630. 개요   631. 개요   632. 개요   633. 개요   634. 개요   635. 개요   636. 개요   637. 개요   638. 개요   639. 개요   640. 개요   641. 개요   642. 개요   643. 개요   644. 개요   645. 개요   646. 개요   647. 개요   648. 개요   649. 개요   650. 개요   651. 개요   652. 개요   653. 개요   654. 개요   655. 개요   656. 개요   657. 개요   658. 개요   659. 개요   660. 개요   661. 개요   662. 개요   663. 개요   664. 개요   665. 개요   666. 개요   667. 개요   668. 개요   669. 개요   670. 개요   671. 개요   672. 개요   673. 개요   674. 개요   675. 개요   676. 개요   677. 개요   678. 개요   679. 개요   680. 개요   681. 개요   682. 개요   683. 개요   684. 개요   685. 개요   686. 개요   687. 개요   688. 개요   689. 개요   690. 개요   691. 개요   692. 개요   693. 개요   694. 개요   695. 개요   696. 개요   697. 개요   698. 개요   699. 개요   700. 개요   701. 개요   702. 개요   703. 개요   704. 개요   705. 개요   706. 개요   707. 개요   708. 개요   709. 개요   710. 개요   711. 개요   712. 개요   713. 개요   714. 개요   715. 개요   716. 개요   717. 개요   718. 개요   719. 개요   720. 개요   721. 개요   722. 개요   723. 개요   724. 개요   725. 개요   726. 개요   727. 개요   728. 개요   729. 개요   730. 개요   731. 개요   732. 개요   733. 개요   734. 개요   735. 개요   736. 개요   737. 개요   738. 개요   739. 개요   740. 개요   741. 개요   742. 개요   743. 개요   744. 개요   745. 개요   746. 개요   747. 개요   748. 개요   749. 개요   750. 개요   751. 개요   752. 개요   753. 개요   754. 개요   755. 개요   756. 개요   757. 개요   758. 개요   759. 개요   760. 개요   761. 개요   762. 개요   763. 개요   764. 개요   765. 개요   766. 개요   767. 개요   768. 개요   769. 개요   770. 개요   771. 개요   772. 개요   773. 개요   774. 개요   775. 개요   776. 개요   777. 개요   778. 개요   779. 개요   780. 개요   781. 개요   782. 개요   783. 개요   784. 개요   785. 개요   786. 개요   787. 개요   788. 개요   789. 개요   790. 개요   791. 개요   792. 개요   793. 개요   794. 개요   795. 개요   796. 개요   797. 개요   798. 개요   799. 개요   800. 개요   801. 개요  </pre>	



## 2) 한국어 문화 수업 자료 만들기를 위한 교재 분석

위의 과정에서 학습 데이터에 등장했거나 그와 비슷한 표현을 찾아냈던 것을 볼 수 있었는데, 눈여겨볼 것은 ‘-하자’ 형태의 지시문이 아닌 ‘-무엇일까?’로 끝나는 지시문을 찾아낸 것이다. 혹시 ‘무엇’과 같은 다른 질문 패턴을 학습하였는지 궁금증이 생겨 다른 교재에도 적용해 보았다. 학습 데이터의 거의 모든 지시문이 ‘-자’ 형태였기 때문에 추출할 수 있는 문장은 없을 것이라 예상했지만, 확인된 결과에 따라 모델을 어떻게 개선하면 좋을지 파악할 수 있을 것 같아 시도해 보았다. <한국문화의 이해> 과목에서 한국어 문화 자료를 만들어야 했고, 이 때 한국어 교재의 문화 파트를 분석하면서 도구를 사용해 보았다.

The image contains two side-by-side screenshots of a terminal window. Both screenshots show the command: `(kobert_final_env) C:\final>python txb_inference.py input_data\sejong1.txt --model_dir instruction_classifier_model`.  
The left screenshot shows the output for 'sejong1.txt':  
모델 로드 중: instruction\_classifier\_model  
모델 로드 성공.  
[ 최종 추출 결과 ]  
총 문장 수: 113  
추출된 지시문 수: 0  
추출된 지시문이 없습니다.  
The right screenshot shows the output for 'sejonghr\_1A.txt':  
모델 로드 중: instruction\_classifier\_model  
모델 로드 성공.  
[ 최종 추출 결과 ]  
총 문장 수: 35  
추출된 지시문 수: 0  
추출된 지시문이 없습니다.

예상대로 추출된 지시문은 없었음을 확인하였다. 한국어 교재 속 지시문을 수집하여 학습 데이터로 활용한다면 한국어 교재에도 활용할 수 있을 것으로 보인다.

## 3) 그 외 활용 가능성

이 도구를 주로 이번 학기에 수강한 강의에서 사용하기는 했으나, 앞으로도 매 학기 전공 수업에서도 계속해서 활용할 수 있을 것으로 예상되며, 성능을 더 개선한다면 실제 교육 현장에서도 유용하게 사용할 수 있을 것으로 예상된다.

## 5. 결론

### 5.1. 배운 점

- 프로젝트에 활용할 모델을 선정하는 과정에서 다양한 모델의 역할과 목적을 알게 되었다. 이전까지는 BERT라는 모델에 대해 기본적인 내용만 알고 있었고, 실제로 사용해 본 적이 없었는데, 이번 프로젝트를 진행하면서 BERT뿐만 아니라 KoBERT라는 모델이 존재한다는 것을 알게 되었고, 실제로 활용해 보면서 모델의 기능을 알게 되었다. 또한, 어떤 모델을 활용할지 찾아보고 선정하는 과정에서 그 밖의 다양한 모델에 대해서도 많이 알게 된 시간이었다.
- 간단한 도구를 만든다고 해도 수집해야 할 데이터의 규모나 형식, 전처리 등 고려해야 할 것이 매우 많고 복잡하다고 느꼈다. 특히, 전처리만 잘 되어도 모델의 성능에 크게 영향을 준다는 것을 깨달았고, 그만큼 그 중요성을 크게 체감할 수 있었다.

### 5.2. 개선 방향

- 1) pdf 파일을 txt파일로 수동으로 변환하는 과정을 자동화한다.

프로젝트를 진행하면서 가장 생각지도 못했던 난관은 바로 모델에 입력하는 자료의 형태였다. 교과서 파일이 일반적인 pdf로 구성되어 있다고 생각하고 프로젝트를 계획하였지만, 모델

학습 과정에서 교과서 파일은 사진을 pdf로 변환한 형태로 제공되거나, e-book 형태로 배포되고 있었다는 것을 깨닫게 되었다. 이 때문에 OCR을 활용하여 직접 수동으로 텍스트를 추출하여 txt 파일로 저장해서 모델에 입력할 수밖에 없었는데, 이 과정 때문에 오히려 시간이 더 많이 소요되기도 했다. 이 도구가 계속해서 사용되려면 이 번거로운 과정에 대한 개선이 꼭 필요하다. OCR 과정을 모델에 넣거나, 그와 비슷한 수준의 자동화 과정이 필요하다.

2) 데이터를 더 정교하게 전처리한다.

실행 결과를 보면, 띠어쓰기의 문제로 출력된 결과가 정돈되지 않은 형태를 보이고 있고, 단원의 제목과 같은 불필요한 텍스트가 함께 출력되고 있다. 띠어쓰기나 문장부호를 처리하고 불필요한 표현은 사전에 없애서 정제된 데이터를 만든다.

3) 학습 데이터의 양과 종류를 확대한다.

학습 데이터의 양을 늘려 국어 교과서에서 한국어 교재나 다른 교재로 확장한다면 더 광범위한 활용이 가능할 것으로 예상된다.