

# 영어분석을위한기계학습 – 6th Assignment Final Report

프로젝트명: BonCahier AI — 수업 자료 자동 요약과 번역 및 핵심 정리 도구

과목명: 영어분석을위한기계학습

학번: 201903774

이름: 한형준

전공: 언어인지과학과

제출일: 2025-12-17

## 1. 프로젝트 개요

현대 대학 수업 환경에서 학생들은 여러 과목의 강의 자료(PPT, PDF)를 동시에 관리해야 하고, 과제·시험 일정이 몰리는 시기에는 짧은 시간 안에 많은 슬라이드를 정리해야 하는 문제가 있다. 특히 다음과 같은 어려움이 존재한다.

- 텍스트 양이 많은 슬라이드(특히 이론 설명 위주 과목)의 경우,  
직접 요약 노트를 만드는 데 시간이 과도하게 소요됨
- 수업 자료가 영어 등 외국어로 작성된 경우,  
번역과 이해를 동시에 해야 하므로 인지적 부담이 큼
- 한 학기 동안 받은 모든 파일(PPT, PDF)을 시험 직전에 다시 열어보며  
핵심만 정리하기가 매우 비효율적

이를 해결하기 위해, 본 프로젝트에서는 강의 자료(PPT, PDF)를 업로드하면 슬라이드/페이지 단위로 텍스트를 추출하고, KoBART 기반 모델로 핵심 요약과 번역을 자동 생성하는 도구 “BonCahier AI”를 설계·구현하였다.

## 2. 진행 과정 요약 (Assignment 3–5)

### 2.1 Assignment 3 – 프로젝트 제안서 작성

Assignment 3에서는 먼저 문제 정의와 서비스 아이디어를 구체화했다.

- 목표 서비스:
  - PPT/PDF 강의 자료 업로드 → 핵심 요약 + 번역 + 간단한 설명 자동 생성
- 주요 기능:

- 파일 업로드(PPT/PDF) 및 텍스트 추출(PyMuPDF, python-pptx)
  - 문단/슬라이드 단위 요약
  - 외국어 자료에 대한 한국어 번역
  - 결과를 Markdown/PDF 형식으로 저장
  - 기대 효과:
    - 시험 대비, 과제 준비, 발표 준비 시  
"핵심 요약 노트"를 자동으로 얻어 **복습 효율 극대화**
- 이 단계에서는 실제 구현보다는 전체 그림(UX 시나리오, 모델 후보, 기술 스택)을 설계하는 데 집중했다.

## 2.2 Assignment 4 – 데이터 수집 및 분석

Assignment 4에서는 강의 슬라이드 요약 태스크와 유사한 구조를 갖는 **요약 데이터셋을 수집·분석(EDA)**하였다.

- 한국어 요약 데이터:
  - daekeun-ml/naver-news-summarization-ko
  - 약 2만 건 이상의 뉴스 기사(document)와 요약문(summary) 쌍
  - 결측치 거의 없고, 카테고리 분포·문장 길이 분포 등을 확인
- 영어 요약 데이터:
  - gopalkalpande/bbc-news-summary
  - BBC 뉴스 기사와 요약문 쌍

EDA 결과, 한국어 네이버 뉴스 데이터는 **요약 태스크의 대표적인 벤치마크로 활용 가능**했고, 문장 길이와 요약 길이 분포도 강의 슬라이드의 "설명형 텍스트"와 어느 정도 유사하다고 판단하였다. 따라서 Assignment 5에서는 이 데이터셋을 **슬라이드 요약 태스크의 proxy corpus**로 사용하기로 했다.

## 2.3 Assignment 5 – KoBART 요약 모델 학습 및 평가

Assignment 5에서는 HuggingFace의 gogamza/kobart-base-v2를 기반으로 네이버 뉴

스 요약 데이터에 대해 KoBART 한국어 요약 모델을 미세조정(fine-tuning) 했다.

- 데이터 전처리 및 분할:
  - Train/Validation/Test = 8: 1: 1 비율
  - 재현성을 위해 seed = 42 고정
- 학습 설정:
  - MAX\_SOURCE\_LENGTH = 512, MAX\_TARGET\_LENGTH = 128
  - BATCH\_SIZE = 4, NUM\_TRAIN\_EPOCHS = 4, LEARNING\_RATE = 3e-5
- 모델 저장:
  - 학습 완료 후 OUTPUT\_DIR에 모델과 토크나이저 저장 → 이후 evaluation.ipynb, inference.ipynb, 그리고 이번 Assignment 6의 boncahier\_service.py에서 재사용 가능하도록 구조화

평가는 ROUGE-1, ROUGE-2, ROUGE-L 지표를 중심으로 이루어졌으며, 기본 KoBART 모델 대비 한국어 뉴스 요약 품질이 향상되는 것을 목표로 하였다.

지표	값
eval_loss	0.505
ROUGE-1	35.62
ROUGE-2	14.46
ROUGE-L	34.99
ROUGE-Lsum	44.90

### 3. Assignment 6 – 모델 서비스화 구조

Assignment 6의 목표는 지금까지 학습한 모델을 실제로 사용할 수 있는 서비스 형태로 만들고, 이를 최소 5회 이상 사용해 본 뒤, 전체 과정을 최종 보고서에 정리하는 것이다.

이번 단계에서 구현한 핵심은 **CLI 기반 요약·번역 도구**인 boncahier\_service.py이다.

#### 3.1 시스템 아키텍처 개요

##### 1. 입력 단계

- 사용자로부터 .pdf 또는 .pptx 파일 경로를 인자로 받는다 (-i/--input).

## 2. 텍스트 추출 단계

- PDF: PyMuPDF(fitz)를 사용하여 각 페이지의 텍스트를 추출
- PPTX: python-pptx를 사용하여 각 슬라이드의 텍스트를 추출

## 3. 언어 감지(간단 heuristic)

- 텍스트 내에 한글(유니코드 범위 \uac00-\ud7a3)이 포함되었는지 확인하여 대략적으로 “한국어 vs 비한국어(주로 영어)”를 구분

## 4. 번역 단계 (옵션)

- 비한국어(주로 영어)로 추정되는 텍스트는 Helsinki-NLP/opus-mt-tc-big-en-ko 모델로 영→한 번역 후 요약
- 한국어 요약 결과는 필요시 Helsinki-NLP/opus-mt-ko-en으로 한→영 번역

## 5. 요약 단계

- KoBART 기반 요약 모델(gogamza/kobart-base-v2 또는 fine-tuned 모델)로 문단 단위 한국어 요약 수행

## 6. 결과 저장

- 각 페이지/슬라이드별로
  - 원문 텍스트
  - 한국어 요약
  - (옵션) 영어 번역 요약
  - 비고(번역 사용 여부, 에러 등)
- 위의 내용을 Markdown 파일(summary\_{파일이름}\_{타임스탬프}.md)로 저장

이 구조를 통해, 하나의 커맨드로 강의 슬라이드 전체에 대한 요약 노트를 얻을 수 있다.

### 3.2 주요 스크립트 설명 – boncahier\_service.py

- extract\_from\_pdf(path)
  - PyMuPDF를 이용해 PDF 각 페이지의 텍스트를 추출하고, {"id": "page\_1", "kind": "page", "text": "..."} 형태 리스트로 반환.

- extract\_from\_pptx(path)
  - python-pptx로 슬라이드별 텍스트를 모아 {"id": "slide\_1", "kind": "slide", "text": "..."} 형태로 반환.
- load\_summarization\_model(model\_dir, model\_name)
  - model\_dir가 존재하면 그 디렉토리(학습 완료된 KoBART)를 사용
  - 없으면 model\_name (기본: gogamza/kobart-base-v2)에서 모델 로드
- summarize\_text(...)
  - KoBART 요약 모델로 한국어 문단 요약 수행
- translate (text, direction, cache)
  - direction이 "en-ko" 또는 "ko-en"인 경우, 적절한 Helsinki-NLP 번역 모델을 로드하여 번역 수행
- process\_file(...)
  - 입력 파일 하나를 처리해서 최종 Markdown 요약 파일을 생성하는 상위 함수
- main ()
  - CLI 인자를 파싱하고, process\_file()을 호출하여 실제로 요약을 수행

#### 4. 실제 사용 결과 (5회 이상)

회차	사용 일 시	입력 파일 이름	파일 유 형	언어	주요 목 표	소감
1	2025-12-06 11:00	201903774한형준15주발표자료.pptx	ppt	한국어	요약	핵심 키워드 중심으로 요약은 나쁘지 않다.
2	2025-12-06 11:50	ml4e-lecture-week13.pdf	pdf	영어	요약&번역	영어->한국어 번역에 심각한

						오류가 있다.
3	2025-12-06 12:03	2025fall_7.DNN과 RNN.pdf	pdf	한국어	요약&번역	한국어->영어 번역에도 다소 오류가 있으나 영어->한국어 번역에 비해 심각하지 않음
4	2025-12-06 12:30	NLP_Lec3&4_NLP_TM_slides.pdf	pdf	영어	요약	번역에 비해 핵심 키워드 중심을 찾지만 전처리 과정이 필요함
5	2025-12-06 12:40	PythonReview.pdf	pdf	영어	요약	프로그래밍 언어가 포함된 슬라이드는 오류가 다수 검출됨

스크린샷:

1회차:

# 4번 퀸 - 강의 PPTX/PDF 업로드 후 실행 (업로드 방식)  
# from google.colab import files  
# uploaded = files.upload() # 여기서 .pdf 또는 .pptx 파일 선택

2019037747번은 15주별자료.pdf  
2019037747번은 15주별자료.pptxapplication/vnd.openxmlformats-officedocument.presentationml.presentation - 1728734 bytes, last modified: 2023. 11. 19. - 100% done  
Saving 2019037747번은 15주별자료.pdf to 2019037747번은 15주별자료.pdf

# 딕셔너리 형태의 파일명 확인  
list(uploaded.keys())

[1]: 2019037747번은 15주별자료.pptx

# 파일 이름 꺽고 아래와 같이 실행하기  
list(uploaded.keys())[0] # 첫 번째 업로드 파일  
input\_path = f'{content[uploaded.keys()[0]]}'  
print(input\_path)  
# python /content/initial/benchmark\_service.py  
# -i {input\_path} # --model\_dir={MODEL\_DIR}  
# --start-page 2  
# --end-page 21  
# --no-translate  
# -o {OUTPUT\_DIR}

... INPUT: /content/2019037747번은 15주별자료.pptx  
[Serialization] Loading model from: /content/drive/MyDrive/benchmark/models/robot\_ko.nms  
[TensorRT] TensorRT version: 8.0.1  
2023-12-06 02:04:01.253189 External /local/xla/stubs/executor/cuda/cuda\_tf.cc(467) Unable to register cufft factory: Attempting to register factory for plugin cufft when one has already been registered  
[TensorRT] The logs messages before absl::InitializeFlag() is called are written to /tmp/STCER  
#0000 00:00:00.764995641 283162 2151 cudaBias.cc(407) Unable to register cuDAS factory: Attempting to register factory for plugin cuDAS when one has already been registered  
#0000 00:00:00.764995641 300752 2151 computation\_iplace.cc(177) computation placer already registered. Please check linkage and avoid linking the same target more than once.  
#0000 00:00:00.764995641 300752 2151 computation\_iplace.cc(177) computation placer already registered. Please check linkage and avoid linking the same target more than once.  
#0000 00:00:00.764995641 300752 2151 computation\_iplace.cc(177) computation placer already registered. Please check linkage and avoid linking the same target more than once.  
#0000 00:00:00.764995641 300802 2151 computation\_iplace.cc(177) computation placer already registered. Please check linkage and avoid linking the same target more than once.  
2023-12-06 02:04:01.306928 [TensorFlow] core/lib/tensorflow/cpu\_feature\_guard.cc(210) This TensorFlow binary is optimized to use available CPU instructions in performance-critical operations.  
[TensorFlow] core/lib/tensorflow/cpu\_feature\_guard.cc(210) To disable this behavior, set the TF\_CUDNN\_DISABLE=1 environment variable.

✓ 외부 링크 확인  
import glob, os  
%matplotlib inline

3회차:



마지막 단계 준비 도구

6thAssignment.ipynb - Colab

Google Drive

한국외국어 대학교 e-Class Sy

HUFS-LAI-MAE-2025-2/assig

(102) Vndeopop is a Bad Person

(102) バーチャルアート

파일 수정 보기 삽입 터미널 도구 도움말

Q. 캡처어 + 코드 + 텍스트 > 모두 실행

```

# ④ 번째 파일 - 각의 PPTX/PDF 업로드 후 실행 (업로드 방식)
1 from google.colab import files
2 uploaded = files.upload() # 여기서 .pdf 또는 .pptx 파일 선택
3
4 uploaded = NLP_Lec384_NLP_TM_slides.pdf
NLP_Lec384_NLP_TM_slides.pdf(application/pdf) - 3021497 bytes, last modified: 2025. 12. 6. - 100% done
Saving NLP_Lec384_NLP_TM_slides.pdf to NLP_Lec384_NLP_TM_slides.pdf
['NLP_Lec384_NLP_TM_slides.pdf']

# ⑤ 번째 파일 이름 풀고 아래와 같이 실행하기
1 # 파일 이름 풀고 아래와 같이 실행하기
2 input_fnames = list(uploaded.keys())[0] # 첫 번째 업로드 파일
3
4 uploaded_fnames = f'/content/{input_fnames}'
5
6 lexton(/content/final/final/bancalier-service.py
7 -i '{uploaded_fnames}' # 파일 경로
8 --model_dir={os.path.join('outputs', 'bancalier', 'tfhub'))}
9 --start_page=15 #
10 --end_page=30 #
11 --no_translate #
12 -o '{output_dir}'
```

-- INPUT: /content/NLP\_Lec384\_NLP\_TM\_slides.pdf  
[Summary] Loading model from /content/drive/MyDrive/bancalier/tfhub/Kobart\_ko\_news  
You are using TensorFlow which is not compatible with cuDNN. Please use cuTorch.  
2025-12-06 00:29:36.986302: E external/local\_xtk/xla/stream\_executor/cuda/cuda\_fft.cc:407] Unable to register cuFFT factory: Attempting to register factory for plugin cuFFT when one has already been registered  
WARNING: All log messages before abel::InitialZer0() is called are written to STDERR  
0000 00:00:1764991777.038954 23957 cuda\_bias.cc:1407] Unable to register cuBLAS factory: Attempting to register factory for plugin cuBLAS when one has already been registered  
0000 00:00:1764991777.038954 23957 computation\_place.cc:1771] computation placer already registered. Please check linkage and avoid linking the same target more than once.  
0000 00:00:1764991777.038954 23957 computation\_place.cc:1771] computation placer already registered. Please check linkage and avoid linking the same target more than once.  
0000 00:00:1764991777.038954 23957 computation\_place.cc:1771] computation placer already registered. Please check linkage and avoid linking the same target more than once.  
0000 00:00:1764991777.038954 23957 computation\_place.cc:1771] computation placer already registered. Please check linkage and avoid linking the same target more than once.  
0000 00:00:1764991777.038954 23957 computation\_place.cc:1771] computation placer already registered. Please check linkage and avoid linking the same target more than once.  
To enable the following instruction, add the --use\_x86\_instructions option to your TensorFlow binary.  
✓ 오류이 메시지입니다. 결과 파일은 /content/drive/MyDrive/bancalier/outputs/summary\_NLP\_Lec384\_NLP\_TM\_slides\_20251206-032943.ad

```

# ⑥ 번째 결과 확인
1 # 결과 확인
2 import glob, os
3
```

변수 탐색

PM 12:30 T4(Python 3)

3°C 환경

오늘 12:30 2025-12-06

마지막 단계 준비 도구

6thAssignment.ipynb - Colab

Google Drive

한국외국어 대학교 e-Class Sy

HUFS-LAI-MAE-2025-2/assig

(102) Vndeopop is a Bad Person

(102) バーチャルアート

파일 수정 보기 삽입 터미널 도구 도움말

Q. 캡처어 + 코드 + 텍스트 > 모두 실행

```

# ⑦ 정하는 파일 열어보기
1 # 정하는 파일 열어보기
2 with open(files[3], "r", encoding='utf-8') as f:
3     print(f.read())[:10000] # 앞부분만 확인
```

--

```

## Unit 11 - page_25 (page)
### 원문 텍스트
...
Who TM?
TM has become more practical thanks to
the development of big data platforms
Deep learning algorithms that can analyze massive sets of
structured data
Mining and analyzing text helps organizations
Find potentially valuable business insights in corporate
documents, news articles, call center logs, verbatim survey
comments, social media posts, medical records and other
sources of text-based data
TM is used in
Increasingly incorporate into AI chatbots and virtual agents
that provide automated responses to customers as part of
their marketing, sales and customer service operations
Nov. 25 & Dec. 10, 2025
Introduction to Language Technology 2025
26
...
### 한국어 요약 결과
Introduction to Language Technology 2025
'to AI chatbots and virtual agents, as part of automated responses to customers and logs, verbatim survey, social media posts, medical records and other service operations, their marketing
...
## Unit 12 - page_27 (page)
### 원문 텍스트
...
By Tim (cont'd)
There are many examples of text-based documents (all
in 'electronic' format -)
E-mails, books, news articles, blog posts, corporate Web
pages, customer surveys, résumés, medical records, technical
papers, incident reports, messages/posts on social networking
```

변수 탐색

PM 12:30 T4(Python 3)

FTSE 코리아 지수

오늘 12:33 2025-12-06

5회차:

# #인자는 파일 열어쓰기  
with open('file4', 'r', encoding='utf-8') as f:  
 print(f.read()[:20000]) # 일본문자 확인

## 한국어 디코딩  
Note because x[0] returns True or False according to the same  
as x == None. If the simplified shows ethat all x[0]  
is true, then it's more likely that x[0] is true.  
---  
# Unit 16 - page\_17 (page)  
## 질문 및 풀이  
---  
§ Tuples are defined with parentheses: (1),(2,(3), for example.  
In essence, you can think of a tuple as a list that is not mutable.  
Thus, indexing works as with list objects.  
§ Examples:  
t=(1,(2,(3  
t[1]=100  
# This won't work because it would change t.  
§ If you know that you aren't going to change your object at all, (then tuple is probably a good choice over list.  
§ The tuple() function will convert its argument (to a tuple, (and list) will convert its argument to a list.  
§ Examples:  
§ t=(1,(2,(3  
y=t[1]  
y[1]=100 #This line execution  
# It will work perfectly.  
tuple([1,2,3])=[1,2,3]  
The basic tuple's SomeString  
tuple  
---  
## 한국어 디코딩  
thus, (indexing works as with part list objects) But assignment fails, (1,(2,(3), for example,)) function will convert its argument to list) will convert its, \_ if it's not (you that you know thou  
---  
# Unit 17 - page\_18 (page)  
## 질문 및 풀이

## 로그 발췌 예시:

1회차:

## Unit 8 — slide 10 (slide)

### 원본 텍스트

...

인간 대상 연구에서 발견된 공격적인 행동과 전전두엽 전부 기능 장애의 연관성: 반사회적인 개개인에게는 회백질과 혈역학 반응의 감소, 배외측 전전두엽 피질의 왼쪽, 안와 전두 피질 오른쪽, 전방 대상 피질 오른쪽에 지역적인 뇌 혈류가 관찰되었다.<sup>63)</sup>

쥐 대상 연구: 내측 전전두엽 피질의 기능을 억제하면 수컷 쥐 간의 공격성이 증가하나 활성화하면 공격성이 감소하고<sup>69)</sup>, 세로토닌이 감소할 수록 공격성이 증가한다.<sup>70)</sup>

전두엽 손상이 목표에 대한 공격적인 반응과 관련이 있고, 측두엽 손상이 격렬하지만 집중되지 않은 분노와 공격성과 관련이 있다는 가정도 있다.<sup>72)</sup>

배외측 전전두엽 피질의 손상이 자기 제어를 실패하게 하여 분노 조절에 결함을 일으키고 도발 당하는 경우 공격적인 충동으로 이어진다.<sup>73)</sup>

공격적인 행동과 전전두엽 피질 장애의 상호 연관성

...

### **### 한국어 요약 결과**

반사회적인 개개인에게는 회백질과 혈역학 반응의 감소, 배외측 전전두엽 피질의 왼쪽, 안와 전두 피질 오른쪽, 전방 대상 피질 오른쪽에 지역적인 뇌 혈류가 관찰되었다.<sup>73)</sup>

공격적인 행동과 전전 두엽 피질 장애의 상호 연관성: 내측 전 전두엽피질의 기능을 억제하면 수컷 쥐 간의 공격성이 증가하나 활성화하면 공격성이 감소하고<sup>69)</sup>, 세로토닌이 감소할 수록 공격성이 증가한다.<sup>70)</sup>

측두엽 손상이 목표에 대한 공격적인 반응과 관련이

2회차:

**## Unit 7 — page\_67 (page)**

### **### 원본 텍스트**

...

### Adding nonlinearities in self-attention

- Note that there are no element-wise nonlinearities in self-attention;
- stacking more self-attention layers just re-averages value vectors
- Easy fix: add a feed-forward network to post-process each output vector
- $m_i = \text{MLP}_{\text{output}}(i)$
- 

$$= W_2 * \text{ReLU}(W_1 \text{output}_i + b_1) + b_2$$

<https://web.stanford.edu/class/cs224n/>

Attention Is All You Need (<https://arxiv.org/pdf/1706.03762.pdf>)

...

### ### 한국어 요약 결과

웹 사이트는 귀하가 웹 사이트를 탐색하는 동안 귀하의 경험을 향상시키기 위해 쿠키를 사용하는데 이 쿠키들 중에서 필요에 따라 분류 된 쿠키는 웹 사이트의 기본적인 기능을 수행하는 데 필수적이므로 브라우저에 저장되며 또한 웹 사이트가 사용 방식을 분석하고 이해하는 데 도움이되는 제 3 자 쿠키를 사용한다. 이러한 쿠키를 거부 할 수도 있다. 그러나 이러한 쿠키 중 일부를 선택 해제하면 검색 환경에 영향을 미칠 수 있다. 이러한 쿠키들을 선택 해제 하면 검색 환경에 영향 미칠 수도 있다. 이런 쿠키들 중 일부 선택 해제된 쿠키와 제 3자 쿠키를 사용할 수 있다. 또한 이러한 쿠키들은 웹 사이

### ### (옵션) 영어 번역 요약

The web site uses cookies to improve your experience while you are browsing for web sites, which are required to perform the basic features of web sites, so that cookies are stored in the browser, as well as using my third cookie, which helps the site to analyze and understand how they are used.

### ### 비고

- ※ 원문이 비한국어로 추정되어 `en→ko` 번역 후 요약을 수행했습니다.

3회차:

### ## Unit 10 — page\_12 (page)

### ### 원본 텍스트

...

#### PREVIEW

v순환신경망(Recurrent Neural Network)과 LSTM (Long Short-Term

Memory)

♣순환신경망은 시간성정보를 활용하여 순차데이터를 처리하는 효과적인 학습모델

♣매우 긴 순차데이터(예, 30단어 이상의 긴 문장)를 처리하는데에는 장기의존성을 잘다

루는 LSTM을 주로 사용(LSTM은 선별 기억 능력을 가짐)

v최근에는 순환신경망을 생성모델로 사용

♣예, CNN과 LSTM이 협력하여 자연영상에 주석을다는 문제를 풀

...

### ### 한국어 요약 결과

'Rurrent Neural Network)과 LSTM (Long Short-Term Memory)

어에우긴순차데이터(예, 30단어이상의긴문장)를처리하는데에는장기의존성을잘드러내고 최근에는순환신경망을생성모델로사용

한데예, CNN과LSTM이협력하여자연영상에주석을다는문제를풀을문제를풀다는문제를 풀을 는가?

CNN은 CNN의협력하여 CNN에게주석다는문제를함바, CNTM이협력

### **### (옵션) 영어 번역 요약**

When it comes to dealing with long-term dependencies and long-term sequences, which are more than 30 words long-term sentences, CNN andLSTM have recently been used to create a model that creates problems that address the problem of commenting on natural images in cooperation with CNN?

4회차:

### **## Unit 9 — page\_24 (page)**

#### **### 원본 텍스트**

...

What is TM? (cont'd)

☞Text Mining vs. Data Mining

→Text Mining

) Dealing with unstructured textual data

- Free-form text
  - » No pre-defined or organized in any way
- } Requiring an extra step
  - The unstructured data has to be organized and structured in a way that allows the data modeling and analytics to occur
- } Concerned with the detection of patterns in NL texts
  - Data Mining
- } Handling structured (highly formatted) data
- } Combining disciplines including statistics, AI and machine learning to apply directly to structured data
- } Concerned with the detection of patterns in database

Nov. 25 & Dec. 01, 2025

Introduction to Language Technology 2025

24

...

### ### 한국어 요약 결과

'Tith unstructured data modeling and analytics to occur

Concerned with the detection of patterns in NL texts

Data Mining disciplines including statistics, AI and machine and apply directly patted daa

'Introd

5호차:

## ## Unit 16 — page\_17 (page)

### ### 원본 텍스트

``

§ Tuples(are(defined(with(parentheses: (1,(2,(3),(for(example.

- In(essence,(you(can(think(of(a tuple as(a list that(is(not(mutable.
- Thus,(indexing(works(as(with list objects.

§ But(assignment(fails:

- t=((1,(2,(3)

t[1]=(100

- This(won't(work(because(it(would(change t.

§ If(you(know(that(you(aren't(going(to(change(your(object(at(all,(

then tuple is(probably(a(good(choice(over list.

§ The(built-in tuple() function(will(convert(its(argument(to(a tuple,(  
and list() will(covert(its(argument(to list.

§ Examples:

§ x=((1,(2,(3)

y=(list(x)

x[1]=(100##(Throws(an(exception.

y[1]=(100##(Works(fine;(changes(y.

tuple(list(x))(==x##(True!

The\$basics\$of\$dicts,\$and\$tuples

**tuple**

---

### ### 한국어 요약 결과

Thus,(indexing(works)as(with part list objects)) But(assignment(fails,  
(1,(2,(3),(for(example,)))) function(will(convert(its(argument(to list)) will(covert, its,\_  
If(won't(wou(that(you(know(thou(

## 5. 배운 점 및 한계, 향후 개선 방향

### 5.1 배운 점

- 데이터 수집-EDA-모델 학습-서비스화까지 하나의 흐름을 직접 경험하면서, "단순한 Colab 실습"이 아니라 실제로 사용 가능한 도구를 만드는 과정을 익혔다.
- 한국어 요약에 특화된 KoBART 모델을 사용해 보면서, **인퍼런스 속도·토큰 길이·Beam Search** 파라미터가 결과 품질과 시간에 어떤 영향을 미치는지 체감할 수 있었다.
- PyMuPDF, python-pptx, transformers, HuggingFace Datasets 등을 함께 사용하면서 실제 ML 프로젝트에서 흔히 쓰이는 라이브러리들의 역할을 이해하게 되었다.

### 5.2 한계

- 학습 데이터가 **뉴스 기사**이기 때문에, 교수님 슬라이드의 레이아웃/요약 스타일과는 약간의 차이가 있다.
- 영어 슬라이드의 경우, en→ko → 한국어 요약 → ko→en 순서로 파이프라인이 구성되어 en→ko 번역이 ko→en 번역보다 오류 가능성이 심각하게 높다.
- 텍스트 추출 품질은 슬라이드의 특징에 따라 달라지므로, 표/이미지로 이루어지거나 수식, 프로그래밍 언어, 웹사이트 링크가 다수 포함된 슬라이드는 요약이 어렵다.

### 5.3 향후 개선 방향

- 강의 슬라이드 전용 데이터(예: 실제 PPT 텍스트 + 요약 노트)를 구축하여

KoBART를 **슬라이드 도메인에 특화된 모델로 재학습해 볼 수 있다.**

- 웹 UI(예: Streamlit, Flask)를 없어서, 명령어에 익숙하지 않은 사용자도 **드래그 앤 드롭으로 사용할 수 있는 서비스**로 확장 가능하다.
- 번역 부분은 HuggingFace의 다른 번역 모델이나, 필요하다면 외부 API(예: DeepL, Google Translate)를 조합해 품질을 끌어올릴 수 있다.
- 더욱 효율적이고 성능 높은 영어 데이터 및 모델을 다시 선정하여 en→ko 번역의 품질을 향상시킬 수 있다.

## 6. 결론

본 프로젝트는 **강의 자료 요약·번역 도구 BonCahier AI**를 기획(Assignment 3), 데이터 수집 및 분석(Assignment 4), KoBART 모델 학습 및 평가(Assignment 5), 그리고 실제 사용 가능한 CLI 서비스로 구현(Assignment 6)하는 전 과정을 다루었다.

그 결과, PDF/PPTX 형식의 수업 자료를 입력으로 받아 **슬라이드/페이지 단위로 핵심 내용을 요약하고 번역해 주는 실질적인 도구**를 완성하였고, 이를 실제 강의 자료에 적용하여 **시험 대비용 요약 노트 자동 생성**이라는 초기 목표를 달성할 수 있었다.