

SKKU SCHOLARSHIP

| WEEKLY PROGRESS MEETING (1)

TEAM 스콜라, SKKULAR

강병준 | 김규진 | 박진아 | 장이준 | 주재현

CONTENTS



1. 프론트엔드 팀

2. 백엔드&AI 팀

TEAM. 스킨라
SKKULAR

Brief Schedule

Based on the roles that were largely divided into two categories, we planned the project schedule by dividing it into AI and app parts. The schedule may vary depending on the situation for more efficient project progress.

Table 2. weekly project schedule

		weekly plan													
		2	3	4	5	6	7	8	9	10	11	12	13	14	15
Overall	Project Proposa	Project Proposa			Proposal Latex	Midterm							Finalterm	Final Latex	
AI					data collection data preprocessing	Rule-Based		Q&A Model (BERT)							
App					UI/UX design Database setting	frontend work		backend work connect frontend and backend				Beta Test Debugging			

(a) : AI (b): Back-end (f) : Front-end

~ week 4
Project Proposal

- week 5**
- 1) Proposal Latex 작성
 - 2) App :
UI/UX design (f)
Database setting (b)
 - 3) AI :
Data collection (a,b)
Data preprocessing (a)



1. 프론트엔드 팀

기획 / UI,UX 구성

SKKULAR 핵심 기능

기능1

맞춤형
장학 검색

기능2

맞춤형
장학 알림

기능3

유사 키워드
장학 안내

기능4

지원 마감일 알림

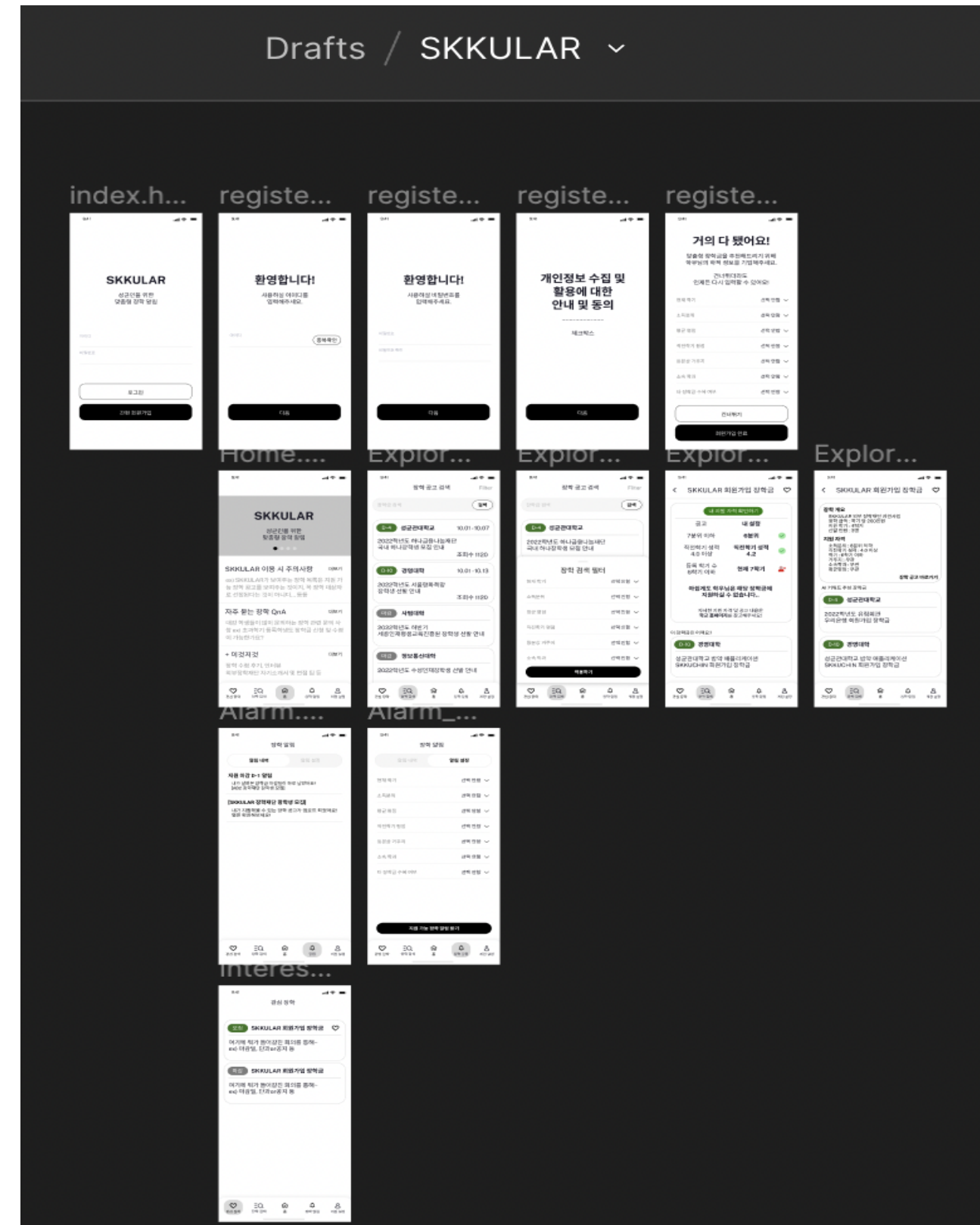
UI/UX 구성 | Figma 사용

1. 프론트엔드 팀 - UI/UX 구성

전체적인 UI/UX 틀

- 1) 로그인 페이지
- 2) 회원가입 페이지
- 3) 홈 화면
- 4) 검색 페이지
- 5) 필터 적용 페이지
- 7) 맞춤형 추천 장학 상세 페이지
- 8) 관심 장학 페이지
- 9) 장학 알림 페이지

+ (10) 계정 설정 페이지)



1. 로그인/회원가입

UI/UX 구성

9:41

SKKULAR

성균인을 위한
맞춤형 장학 알림

아이디

비밀번호

로그인

간편 회원가입

9:41

환영합니다!

사용하실 아이디를
입력해주세요.

아이디

중복확인

다음

9:41

환영합니다!

사용하실 비밀번호를
입력해주세요.

비밀번호

비밀번호 확인

다음

9:41

개인정보 수집 및
활용에 대한
안내 및 동의

~~~~~

체크박스

다음

9:41

거의 다 됐어요!

맞춤형 장학금을 추천해드리기 위해  
학우님의 학적 정보를 기입해주세요.

건너뛰더라도  
언제든 다시 입력할 수 있어요!

현재 학기

선택 안함

소득분위

선택 안함

평균 평점

선택 안함

직전학기 평점

선택 안함

등본상 거주지

선택 안함

소속 학과

선택 안함

타 장학금 수혜 여부

선택 안함

건너뛰기

회원가입 완료

로그인 혹은  
회원가입 페이지로 이동

login.html

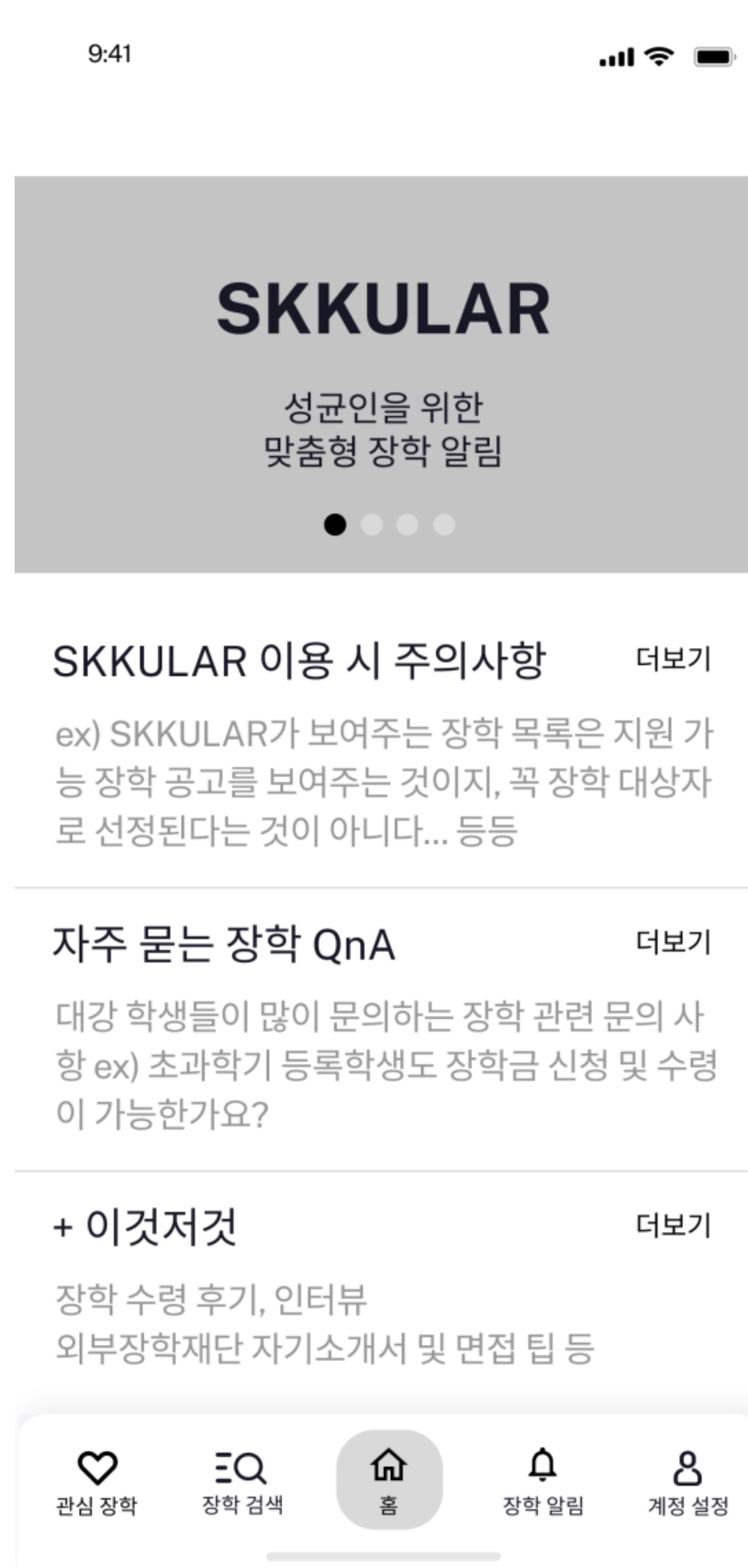
아이디 중복확인 / 비밀번호 일치확인 / 개인정보 수집 동의 / 초기 필터 설정

register1.html ~ register4.html

## 2. 홈 / 검색

home.html

주의사항 및  
Q&A 표시



search.html

전체 장학공지  
리스트 표시



# UI/UX 구성

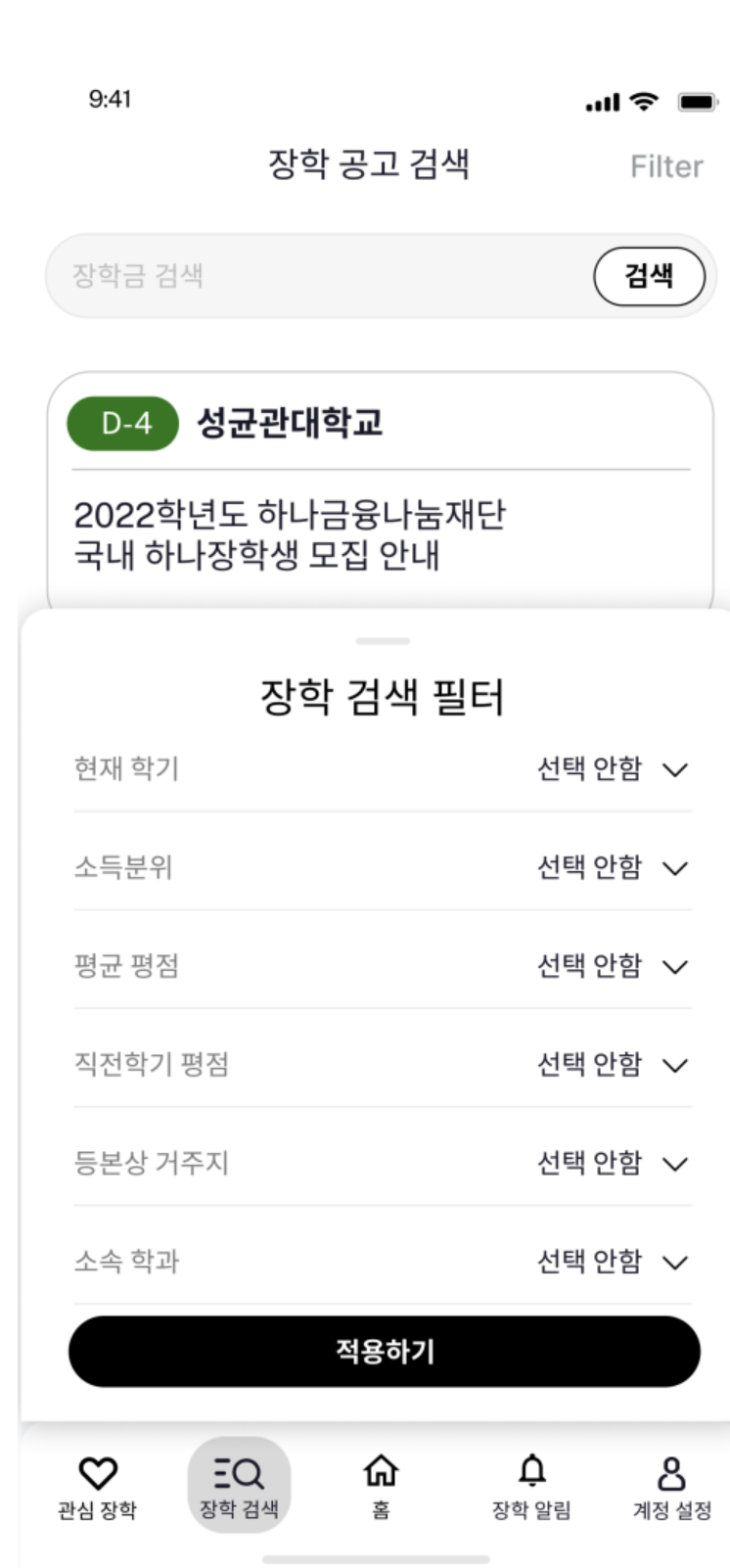


### 3. 필터 / 맞춤형 장학 추천 상세

## UI/UX 구성

filter.html

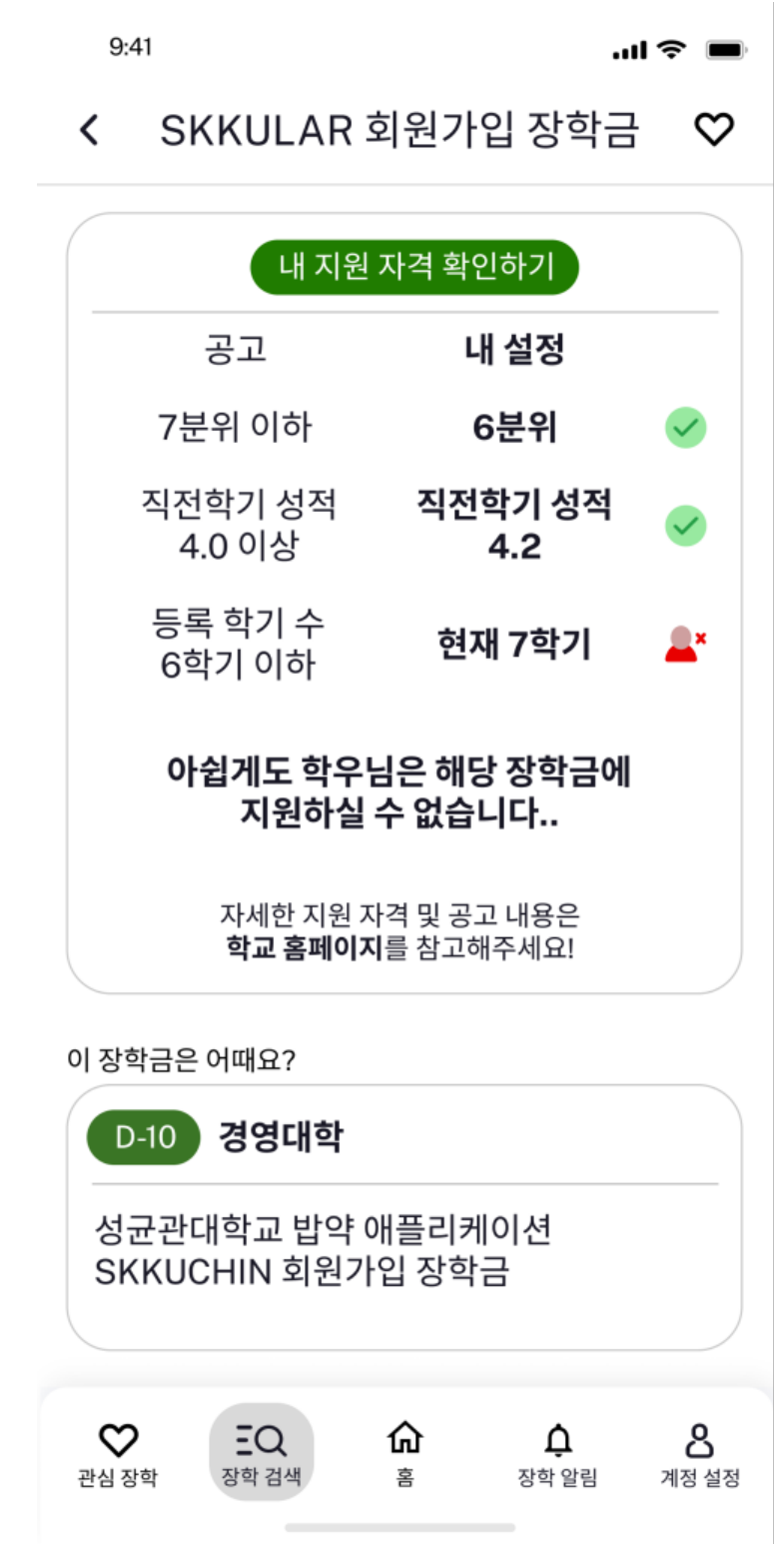
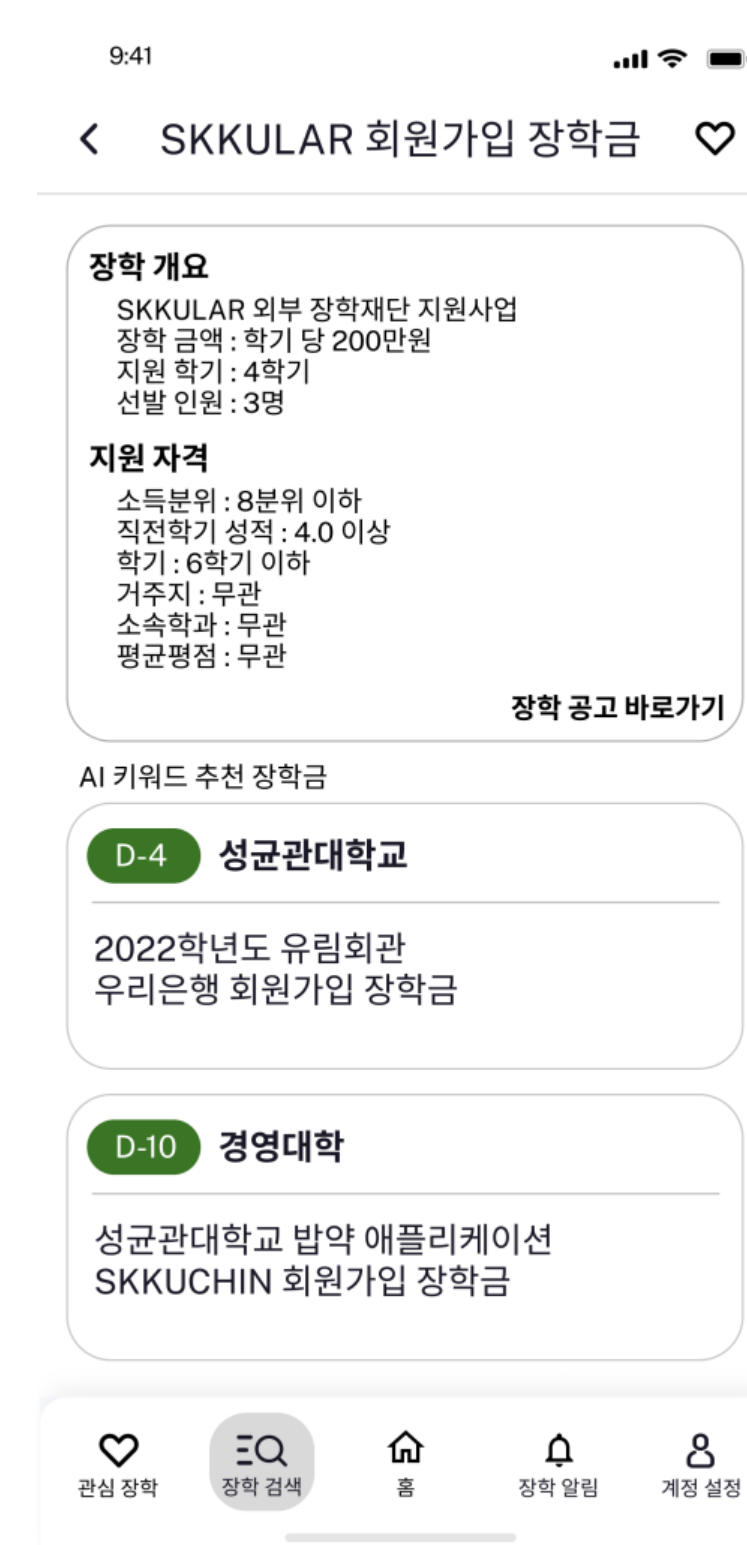
개인 맞춤형  
필터 적용



detail.html

맞춤형 장학 추천,  
내용 상세 보기

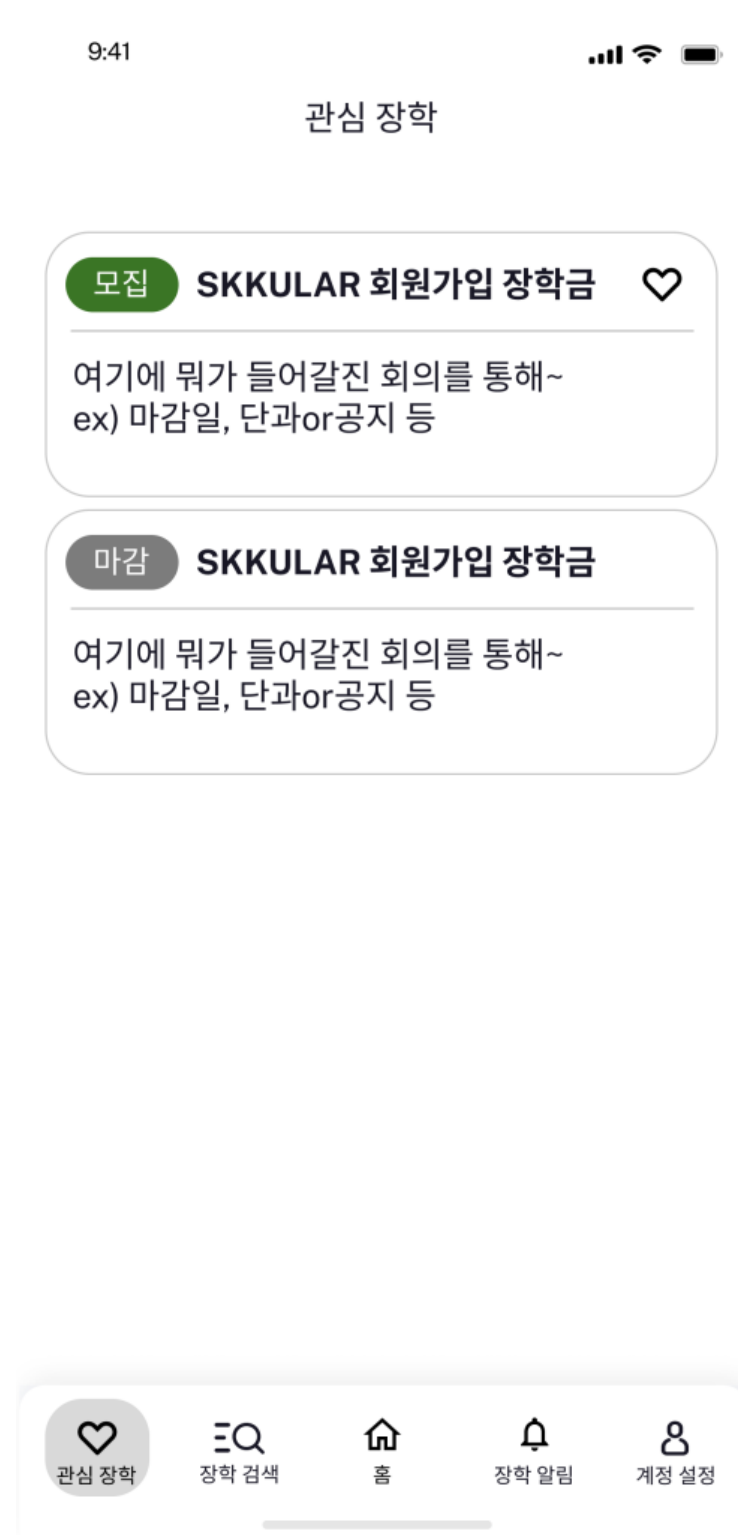
with  
유사 장학금 추천



## 4. 관심 장학 / 장학 알림

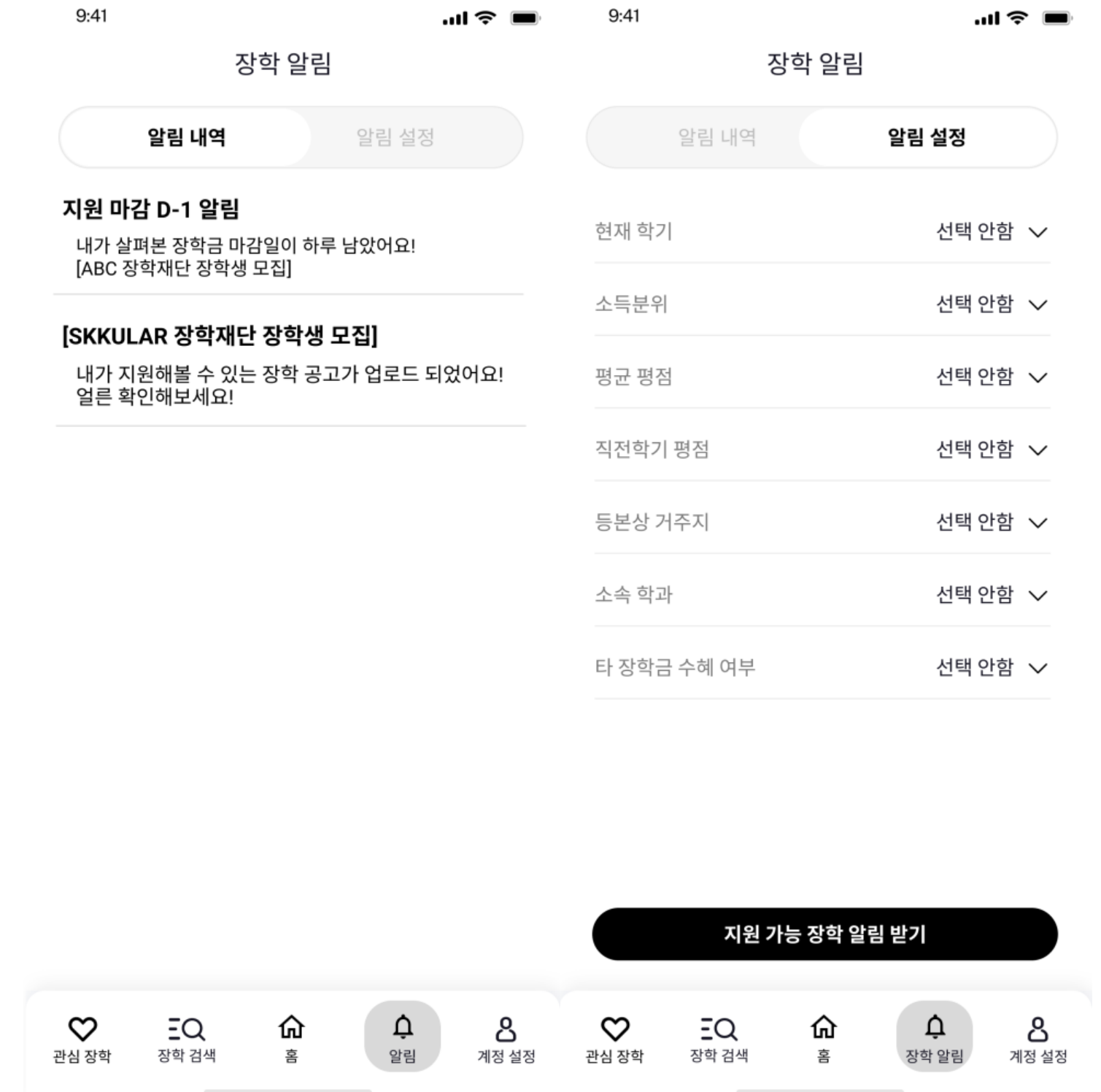
interested.  
html

관심 장학  
스크랩 조회



alarm.html  
alarm\_setting.html

알림 내역 표시,  
알림 on/off 설정



UI/UX 구성



## 2. 백엔드 & AI 팀

### 크롤링

---



# 1. 크롤링 내용 \_ 공지사항 게시판 목록

Ex) 0번째 게시물부터 64번째 게시물 크롤링하고 싶다면,

```
for page_num in crawling_number:
    # 크롤링 우회
    url="https://www.skku.edu/skku/campus/skk_comm/notice06.do?mode=list&&articleLimit=10&article.offset="+str(page_num)
    headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/83.0.4103.61 Safari/537.36'}
    html = requests.get(url, headers = headers)
    result = BeautifulSoup(html.content, "html.parser")

    post_numbers=[]
    # 수집해야 하는 데이터 목록:
    # 1) number/고유번호: No.1806(중복 거르기용)

    post_urls=result.select("#jwx_main_content > div > div > div.container > div.board-name-list.board-wrap > ul > li > dl > dt > a")
    post_numbers=result.select("#jwx_main_content > div > div > div.container > div.board-name-list.board-wrap > ul > li > dl > dd > ul > li:nth-child(1)")
```

1. crawling\_number list에 [0,10,20,30,40,50,60]저장
2. For 문 이용하여 각 page\_num 으로 들어가서 장학금 홈페이지 여러 페이지에 access
3. BeautifulSoup 이용하여 post\_urls와 post\_numbers에 같은 페이지에 있는 여러 게시물의 url과 게시물 번호를 저장

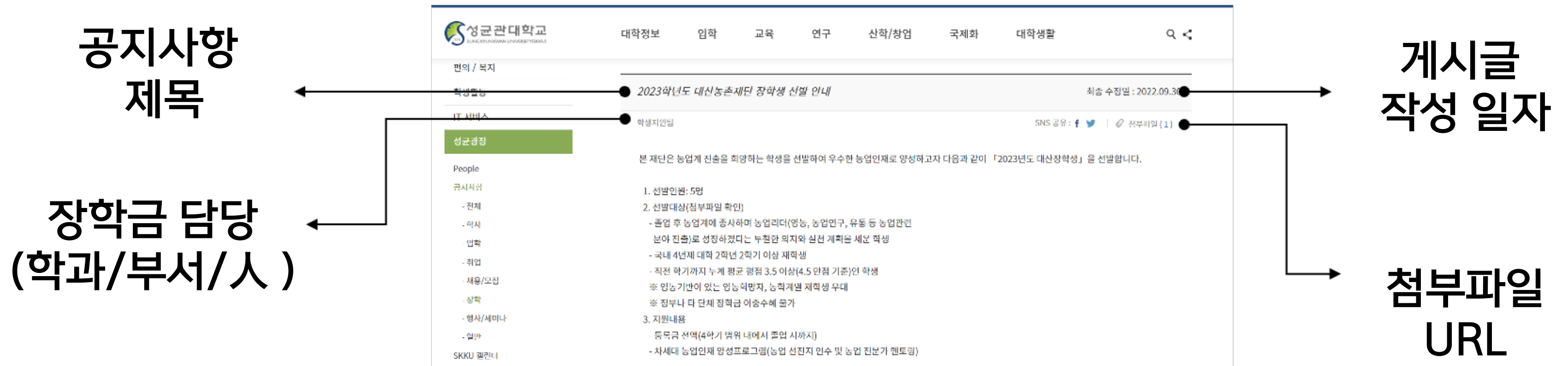
# 1. 크롤링 내용 \_ 공지사항 게시판 목록

```
for i, post_url in enumerate(post_urls):  
  
    print(str(i)+"번째 게시물입니다.")  
    post_url=post_url["href"].replace("&","")  
    post_url="https://www.skku.edu/skku/campus/skk_comm/notice06.do"+post_url  
    headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/83.0.4103.61 Safari/537.36'}  
    html = requests.get(post_url, headers = headers)  
    result = BeautifulSoup(html.content, "html.parser")
```

4. For 문을 이용하여 각 게시물의 post\_url을 얻기

5. BeautifulSoup을 이용하여 각 게시물의 url로 access 하기

## 2. 크롤링 내용 \_ 세부 장학 게시물



## 2. 크롤링 내용 \_ 세부 장학 게시물

```
#1) 홈페이지 고유 number If using all scalar values, you must pass an index
number=post_numbers[i].text

# 세부 페이지 접속: 수집해야 하는 데이터 목록
#2) date/작성일짜: 2022-09-20
date=result.select_one("#jwx_main_content > div > div > div > div > table.board_view > thead > tr > th > span").text.replace("최종 수정일 : ","")
#3) title: 게시물 제목: 2022학년도 하나금융나눔재단 국내 하나장학생 모집 안내
title=result.select_one("#jwx_main_content > div > div > div > div > table.board_view > thead > tr > th > em").text

#4) 게시물 내용
content=result.select_one("#jwx_main_content > div > div > div > div > table.board_view > tbody > tr > td > dl > dd > pre").text
```

1. 홈페이지 고유 번호는 post\_numbers list에 있는 번호 저장
2. 날짜, 제목, 및 게시물 내용은 각각 date, title, content에서 그 부분에 있는 text 크롤링



## 2. 크롤링 내용 \_ 세부 장학 게시물

```
attachments=result.select("#jwxe_main_content > div > div > div > div > table.board_view > tbody > tr > td > div > div.fir > ul > li:nth-child(2) > div > form >
attachment_content=""
attachment_url=""
try:
    for attach in attachments:
        # attach 제목에 attachments
        attach_title=attach.text
        if ("공고문" in attach_title) or ("안내" in attach_title):
            if not os.path.exists("/content/data/"):
                os.makedirs("/content/data/")
            else:
                try:
                    os.remove("/content/data/report_temp.pdf")
                except:
                    pass
            attachment_url="https://www.skku.edu/skku/campus/skk_comm/notice06.do"+attach["href"]
            print(1)
            wget.download(attachment_url,out="/content/data/report_temp.pdf")
            print(2)
            attachment_content= textractor("/content/data/report_temp.pdf")
            print(3)
            os.remove("/content/data/report_temp.pdf")
```

1. 첨부 파일에 있는 url attachments list에 저장 > attach\_title에 첨부파일 제목 저장
2. Attach\_title에서 “공고문“ or “안내” 가 있으면 pdf 다운로드
3. Textract 이용하여 다운로드된 pdf의 내용 저장
4. PDF 삭제

## 2. 크롤링 내용 \_ 세부 장학 게시물

```
if crawling_start==False:

    post_temp=pd.DataFrame({"number":[number], "date":[date], "title":[title],"content":[content], "attachment_url":[attachment_url], "attachment_content":[attachment_
    crawling_start=True
else:
    post_t=post_temp.copy()
    #post_Series([number, date, title,content, attachment_
    utemp=post_t.append([number, date, title,content, attachment_url, attachment_content, current_url, department
    post_temp = post_t.append(pd.r1, attachment_content, current_url, department], index=post_temp.columns), ignore_index=True)

display(post_temp)
```

```
[ ] 1 ##### 구글 드라이브에 지정된 경로로 저장#####
2 df_temp.to_csv("/content/drive/SharedDrives/skkular/data/성균관대학교_전체_장학금.csv",index=False, encoding="utf-8")
```

```
[ ] 1 df_temp
```

1. Pandas dataframe을 만들어서 크롤링한 값을 저장하거나 dataframe에 추가하기
2. 크롤링 된 결과를 저장한 dataframe을 csv file type으로 저장

# 3. 크롤링 결과

## CSV 파일로 저장

| number  | date       | title          | content                                                                                                          | attachment_url                                                                                    | attachment_content | current_url                                               | department |
|---------|------------|----------------|------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------|--------------------|-----------------------------------------------------------|------------|
| No.1805 | 2022.09.20 | 2022학년도 2학기 DE | 당 재단은 DB그룹(구 동부그룹) 설립자 김준기 회장이 기금 출연하                                                                            | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1804 | 2022.09.16 | 2022학년도 (재)광산  | (재)광산장학회에서는 우리 지역의 우수한 인재를 발굴·육성하고 학                                                                             | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1803 | 2022.09.16 | 2022학년도 2학기 인  | 인산인재육성재단에서는 교육비 부담 경감을 통한 지역 인재육성 및                                                                              | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1802 | 2022.09.16 | 2022학년도 2학기 고  | 고양시에서는 대학생 <a href="https://www.skku.edu/">https://www.skku.edu/</a> 본문내용 입력 (2022                               | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1801 | 2022.09.15 | 2022학년도 2학기 63 | 2022년도 2학기 63장학금 장학생을 선발하고자 하오니 아래를 참고                                                                           | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1800 | 2022.09.14 | 2022학년도 하반기 (  | 1. 재단 안내 ㅇ 명 <a href="https://www.skku.edu/">https://www.skku.edu/</a> 1111 재단법인 은평구                             | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1799 | 2022.09.14 | 2022학년도 하반기 (  | 가. 선발인원: 26명(성 <a href="https://www.skku.edu/skku/campus/skku.com">https://www.skku.edu/skku/campus/skku.com</a> | <a href="https://www.skku.edu/skku/campus/skku.com">https://www.skku.edu/skku/campus/skku.com</a> |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1798 | 2022.09.06 | 2022학년도 수험재단   | 수험재단에서 설립자인 동교 김희수 전 이사장의 교육이념을 받들어                                                                              | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1797 | 2022.09.06 | 2022학년도 2학기 지  | 본 재단은 국내 디자인 예술 산업 발전에 밑거름이 될 인재육성을 위                                                                            | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1796 | 2022.09.06 | 2022학년도 일문과학   | 1.선발분야 및 규모: 0 명 내외 가. 전공 제한 없이 지원 가능 나                                                                          | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1795 | 2022.09.06 | 2022학년도 서울장학   | 서울장학재단에서는 <a href="https://www.skku.edu/">https://www.skku.edu/</a> 2022년 <서울희망 디                                | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1794 | 2022.09.06 | 2022학년도 광주광역   | 1. 지원자격 : 공고일 <a href="https://www.skku.edu/">https://www.skku.edu/</a> 재 (재)광주광역시 서                             | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1793 | 2022.09.05 | 2022학년도 하반기    | ■ 신청기간: 2022. 9. 26.(월) ~ 9. 30.(금) (5일간)■ 접 수 처: 애항장                                                            | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1792 | 2022.09.02 | 2022학년도 제 25기  | 「두들장학재단」은 삼성그룹의 창업주故이병철 선대회장을 내조하                                                                                | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1791 | 2022.09.02 | 2022학년도 경주시장   | 1. 선발 인원 및 신청자격□ 선발인원 - 총 450명(일반전형: 423명 특                                                                      | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1790 | 2022.09.02 | 2022학년도 포항시장   | 1. 신청기간 : 2022. 08. 24(수) 09:00 ~ 2022. 09. 20(화) 18:00.                                                         | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1789 | 2022.09.02 | 2022학년도 인송문화   | '仁松文化財團' 이 학업에 뜻이 있고 바른 사회관을 가진 대학생으로                                                                            | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |
| No.1788 | 2022.08.31 | 2022학년도 2학기 도  | 도사랑학회는 1995년 설립되어 매년 경제적 어려움을 겪는 학생들                                                                             | <a href="https://www.skku.edu/">https://www.skku.edu/</a>                                         |                    | <a href="https://www.skku.edu/">https://www.skku.edu/</a> | 학생지원팀      |

- Number : 공지사항 번호
- Date : 공지사항 작성 날짜
- Title : 공지사항 제목
- Content : 공지사항 내용
- Attachment\_url : 첨부파일 url
- Attachment\_content : 첨부파일 내용
- Current\_url : 공지사항 url
- Department : 공지사항 부서

학교 공식 홈페이지 공지사항은 약 1,000개, 단과대 홈페이지 공지사항은 대략 30개씩 크롤링 진행

## 4. 크롤링 역할 분담

### 장이준

- 크롤링 전체적 구현
- 학교 공식 장학 공지,  
학부대학, 유학대학,  
문과대학, 사회과학대학  
장학 공지

### 김규진

- Textract pdf 조사
- 경제대학, 경영대학, 사  
범대학, 예술대학,  
자연과학대학 장학 공지

### 강병준

- Javascript으로 인한  
오류 수정
- 정보통신대학, 글용,  
공과대학, 생명과학대학,  
스포츠과학대학,  
성균융합원 장학 공지

# Trial and error \_ textract pdf 저장



```
5 !wget -N https://github.com/neuml/txtai/releases/download/v3.5.0/tests.tar.gz
6 !tar -xvzf tests.tar.gz
7
8 # Install NLTK
9 import nltk
10 nltk.download('punkt')

1 %%capture
2
3 from txtai.pipeline import Textractor
4
5 # Create textractor model
6 textractor = Textractor()

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

1 pdf_path = "/content/drive/Shared drives/skkular/data/a.pdf"

1 textractor("/content/drive/Shared drives/skkular/data/a.pdf")
```

장학금 공지사항 내 첨부파일이 '신청서' 이면 크롤링 해서는 안 됨!



첨부파일 제목에 “공고문” or “안내” 라는 단어가 있으면 textract 이용하여 pdf 내용 저장



# Trial and error \_ 전체 공지 vs 단과대 공지

```

<thead>...</thead>
<tbody>
  <tr>
    <td>
      <div class="boardView_txtWrap">...</div>
      <dl class="board-write-box board-write-box-v03">
        <dt class="hide replyNone">게시글 내용</dt>
        <dd>
          <pre class="pre"> == $0
          "교육부 국립국제교육원의 해외 석박사 학위과정 지원 장학사
          업을 첨부와 같이 안내합니다. ** 주의 사항 ** 1. 필수 제출
          서류 중 '응시원서'의 총장 직인은 졸업한 단과대학의 행정실
          에 요청하시기 바랍니다. 2. 본 장학사업에 대한 문의는 국립
          국제교육원 국비유학담당자에게 연락하시기 바랍니다. 02-
          3668-1374/1375 kscholar@korea.kr 이상입니다. 국제처 국제
          교류팀"
          </pre>
        </dd>
      </dl>
    </td>
  </tr>
</tbody>

```

전체 공지사항

```

<span style="font-family: Dotum;">
  <span style="font-family: Dotum;">
    <p>
      <span style="font-family: 'Dotum'; font-size: 10pt;">경제
      학과 60학번 동기회에서 경제학과 및 글로벌경제학과 학생들의
      학습의욕 고취와 </span>
    </p>
    <p>
      <span style="font-family: 'Dotum'; font-size: 10pt;">원활
      한 학업수행을 위하여 다음과 같이 장학생을 선발하고자 합니다.
      </span>
    </p>
    <p>
      <span style="font-family: 'Dotum'; font-size: 10pt;">여러
      분들의 많은 신청 바랍니다.</span>
    </p>
    <p>
      <span style="font-family: 'Dotum'; font-size: 10pt;">
      &nbsp;   </span>
    </p>
  </span>

```

단과대 공지사항

전체 공지사항은 <dd>에 있는 text인 반면, 단과대 공지사항은 <p>에 있는 text를 추출

# Trial and error \_ 단과대 공지

```
<ul class="board-view-file-wrap">...</ul>
<div class="board-view-content-wrap board-view-txt">
  <div class="fr-view"> == $0
    <p class="0" style=" color: rgb(0, 0, 0); font-family: "Noto Sans KR", sans-serif; font-size: 17px; text-align: start; background-color: rgb(255, 255, 255);">
      <span style=" font-family: Arial, Helvetica, sans-serif;">대학원팀에서는 2023학년도 1학기 新대학원우수장학생을 다음과 같이 선발할 예정임을 안내하였사오니 우수한 학생 여러분의 많은 지원을 바랍니다.</span>
    </p>
    <p class="0" style=" color: rgb(0, 0, 0); font-family: "Noto Sans KR", sans-serif; font-size: 17px; text-align: start; background-color: rgb(255, 255, 255);">...</p>
    <p class="0" style=" color: rgb(0, 0, 0); font-family: "Noto Sans KR", sans-serif; font-size: 17px; text-align: start; background-color: rgb(255, 255, 255);">...</p>
  </div>
```

```
date=result.select_one("div.board-name-view.board-wrap > div > div.board-view-title-wrap > ul > li:nth-child(3)").text
#date='2022-06-20'
print(date)
#3) title: 게시물 제목: 2022학년도 하나금융나눔재단 국내 하나장학생 모집 안내
title=result.select_one("div.board-name-view.board-wrap > div > div.board-view-title-wrap > h4").text.replace("[경영대학] ", "")

print(title)
#4) 게시물 내용 : p 태그 전체를 감싸고 있는 상위태그를 골라주셔야 합니다!!!!!!!!!!!!
##jwxe_main_content > div > div > div > div > div > div.board-view-content-wrap.board-view-txt
content=result.select_one("div.fr-view").text.replace("제목없음 ", "")
content = content.replace("[경영대학]", "")
print(content)
#5) [첨부파일 제목, 첨부파일 링크]
attachments=result.select("#jwxe_main_content > div > div > div > div > div > ul > li > a")
##jwxe_main_content > div > div > div > div > div > ul > li > a
#####
```

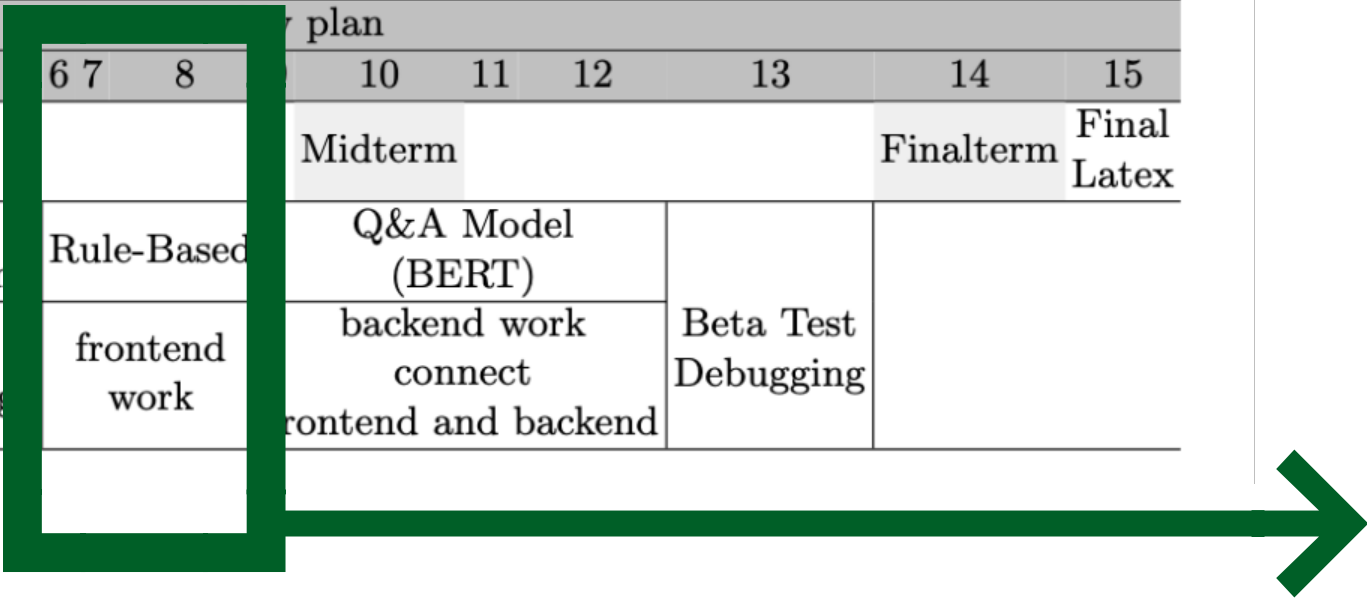
경영대학 등 일부 대학은 javascript에서 selector 값이 보이는 것과 다르게 나오기 때문에  
html에서 값을 직접 select 하는 방식을 취함

# Brief Schedule [next week]

Based on the roles that were largely divided into two categories, we planned the project schedule by dividing it into AI and app parts. The schedule may vary depending on the situation for more efficient project progress.

Table 2. weekly project schedule

|         |                  |   |                                       | Project plan |               |   |  |   |   |   |                                              |    |                        |    |    |  |           |  |             |  |
|---------|------------------|---|---------------------------------------|--------------|---------------|---|--|---|---|---|----------------------------------------------|----|------------------------|----|----|--|-----------|--|-------------|--|
| 2       |                  | 3 | 4                                     | 5            |               | 6 |  | 7 | 8 | 9 |                                              | 10 | 11                     | 12 | 13 |  | 14        |  | 15          |  |
| Overall | Project Proposal |   | Proposal Latex                        |              |               |   |  |   |   |   | Midterm                                      |    |                        |    |    |  | Finalterm |  | Final Latex |  |
|         |                  |   | data collection<br>data preprocessing |              | Rule-Based    |   |  |   |   |   | Q&A Model (BERT)                             |    |                        |    |    |  |           |  |             |  |
| AI      |                  |   |                                       |              |               |   |  |   |   |   |                                              |    |                        |    |    |  |           |  |             |  |
| App     |                  |   | UI/UX design<br>Database settings     |              | frontend work |   |  |   |   |   | backend work<br>connect frontend and backend |    | Beta Test<br>Debugging |    |    |  |           |  |             |  |



(a) : AI (b): Back-end (f) : Front-end

## week 6

### 1) APP :

Database setting (b)

UI/UX design 보완 (f)

Frontend 개발 시작 (f)

### 2) AI :

Data preprocessing (a)

Data labeling (a)

Rule-Based 구현 시작 (a)



---

# Q & A

지금까지 TEAM. SKKULAR 었습니다.  
감사합니다.