

한남대학교 빅데이터

어디 지역 (=분야) 출신이니?

뉴스기사분류

Contents

01 개발배경

02 개발 개요

03 데이터 수집

04 데이터 전처리

05 모델 학습

06 서비스 개발

07 활용방안 및 기대효과

A blurred office scene with a person working at a desk. The desk is equipped with multiple computer monitors, a desk lamp, and various office supplies. The background shows a bright window and other office furniture. The text '개발배경' is overlaid in the center.

개발배경

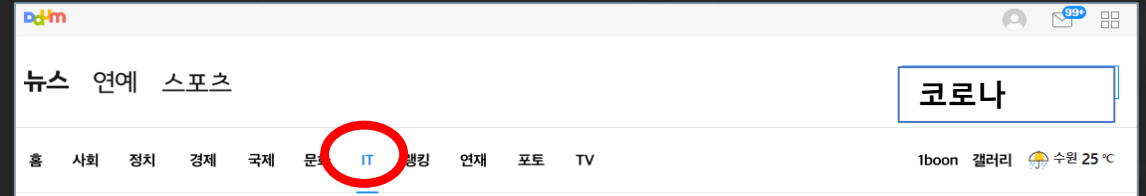
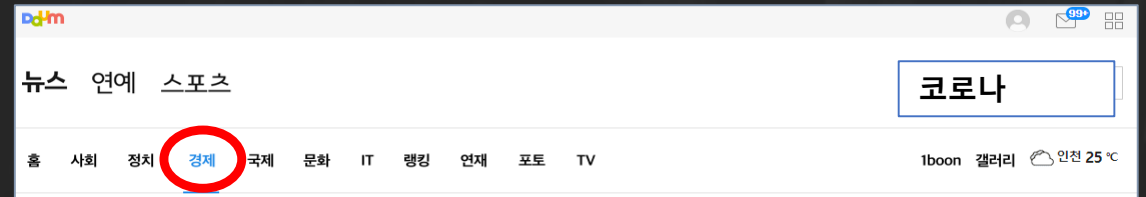
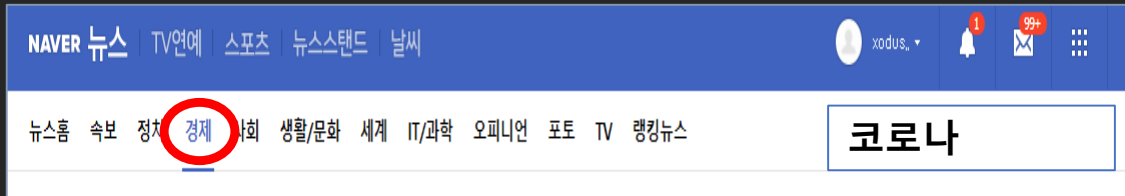
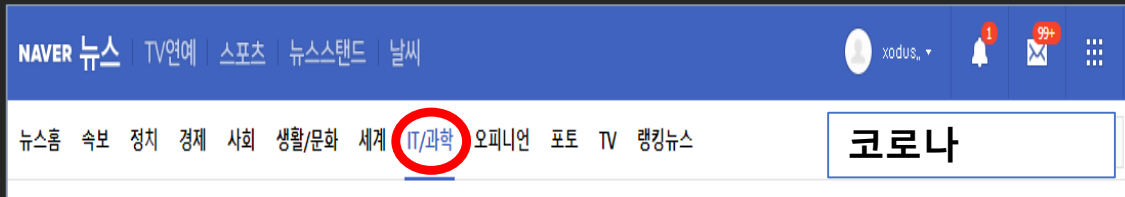
1. 개발 배경 – 주제 선정 동기

방대한 데이터를 찾아보는데 있어서 신속성과 정확성을 갖음과 동시에 원하는 데이터를 찾을 수 있다는 장점을 추가하기 위해서



1. 개발 배경 – 주제 선정 동기

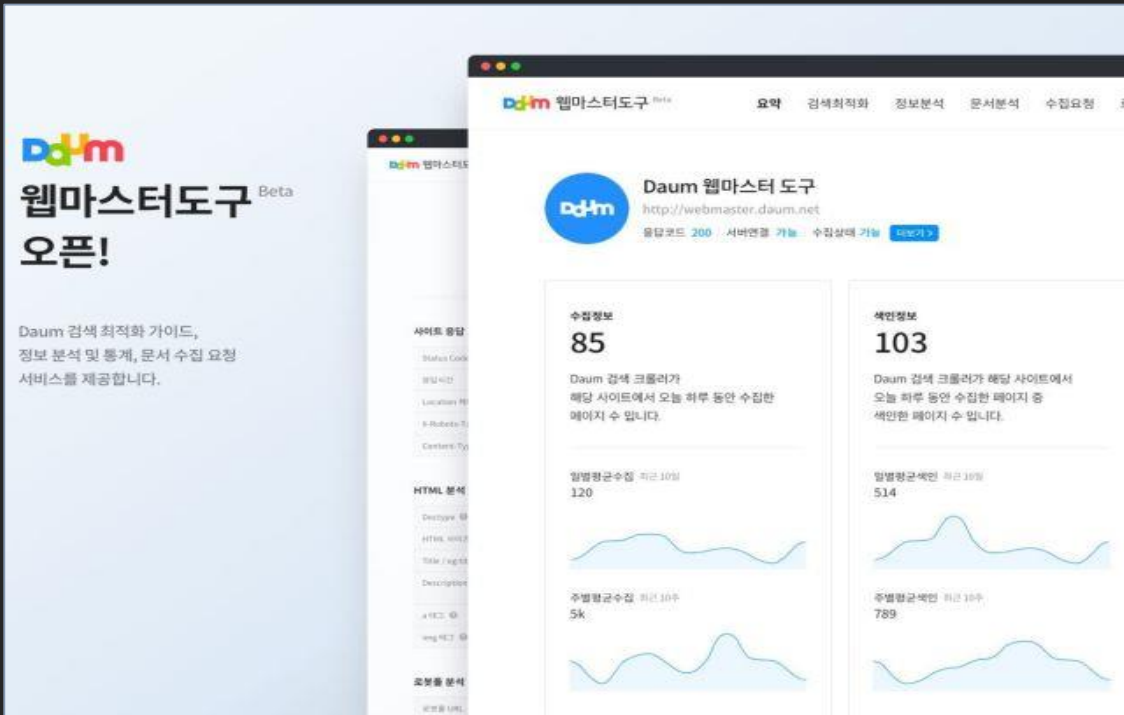
- 대표 포털사이트의 뉴스사이트에서 특정 단어를 검색했을 때 카테고리 별로 검색되어 나오지 않음



1. 개발 배경 – 주제 선정 동기

- 검색 속도 최적화를 돕는 ‘다음 웹마스터 도구 베타’ 공개

카카오, Daum은 검색 엔진 기술력을 활용해 검색 서비스를 향상시키기 위해 ‘Daum 웹마스터 도구 beta’ 를 오픈하여 검색 선순화 활성화로 콘텐츠 공급자와 이용자 모두의 만족도를 높이는데 큰 역할을 할 것으로 기대하고 있다.



네이버, '실검' 일부 개편...내일부터 '연령별 맞춤' 순위 표출

입력 2019.10.30 14:16 | 수정 2019.10.30 14:16

네이버, 실검 첫 화면 '연령별 맞춤' 순위로 변경

실시간 인기

네이버 "집중도 분산 효과 기대"

1 삼성바이오로직스 위탁

플랫폼 전면 개편 나선 카카오에 비해 소극적이라는 평

2 한중일 배터리 삼국지, 1

• 검색 서비스의 불편사항을 개선하는 네이버

네이버 관계자는 오래된 실시간 검색 방식을 개편하여 모든 사람이 같은 실시간 검색어를 보는 **현재의 집중도를 분산시키는 효과**를 볼 수 있을 것으로 기대하고 있다.

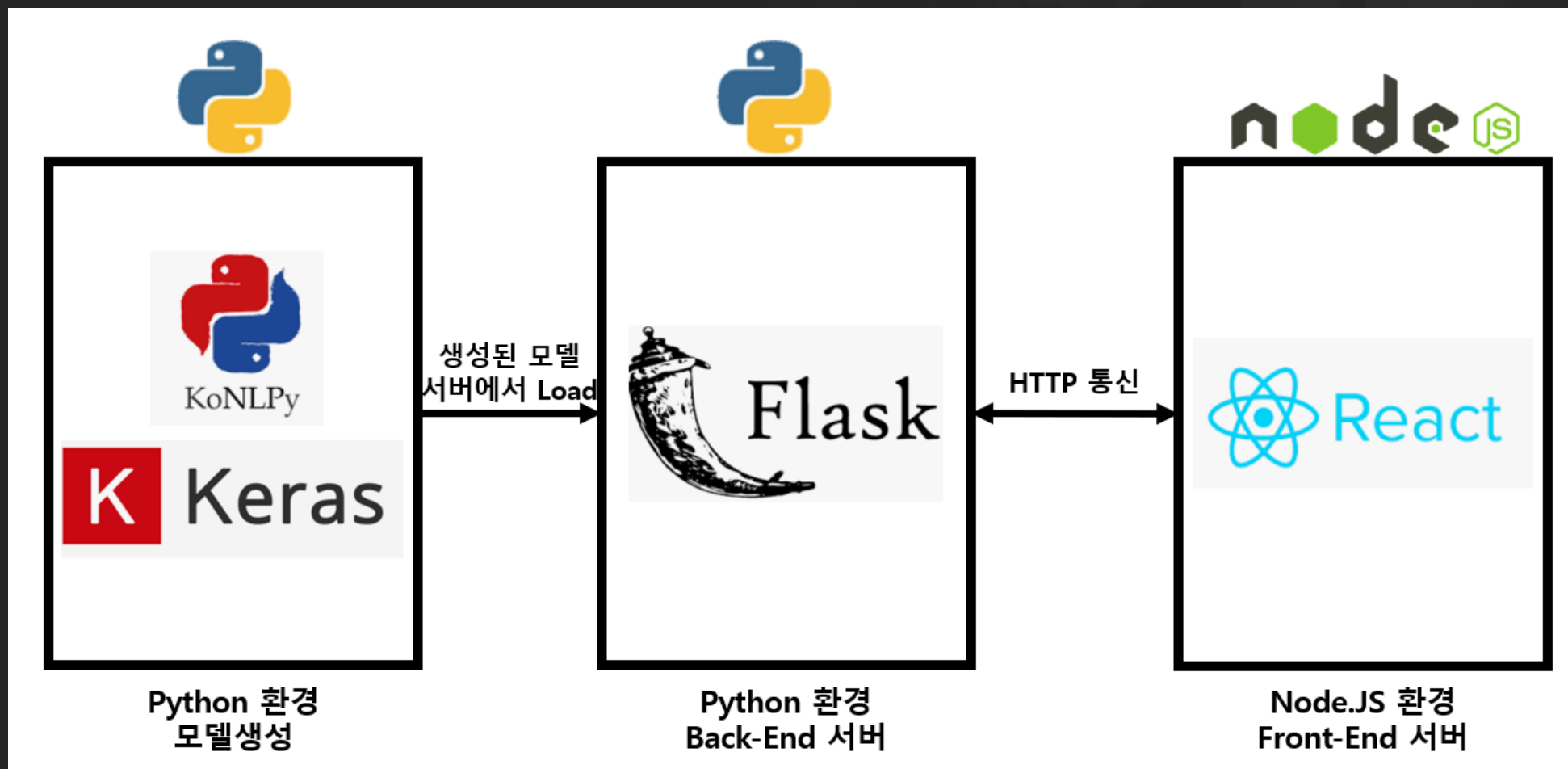
이는 소비자들의 불편, 불만 사항인 실시간 검색 조작의 가능성을 배제하도록 하기 위한 소통의 결과물이라고 볼 수 있다.

출처 : 한국경제



개발개요

02. 개발개요 – 사용 라이브러리



02. 개발개요 - 프로젝트 일정

	21	22	23	24	25	26	27	28	29	30	1	2	3
Data collection													
Data preprocessing													
Professor Counseling													
Data modeling													
Service development													
Document													

A blurred office scene with desks, computers, and a person working in the background. The image is in grayscale and has a soft, out-of-focus quality. In the foreground, there's a desk with a large Apple iMac, a water bottle, and some small potted plants. In the background, a person is sitting at another desk, working on a laptop. The text "데이터 수집" is overlaid in the center.

데이터 수집

03. 데이터 수집- 뉴스 크롤링

원하는 뉴스페이지를 찾아서 URL 복사



뉴스기사 파일 생성

2020-09-02.csv	2020-09-02 오후 4:27	Microsoft Excel ...	28KB
it_2_.csv	2020-09-01 오후 3:53	Microsoft Excel ...	15,545KB
건강_2_.csv	2020-09-01 오후 5:01	Microsoft Excel ...	15,935KB
경제_2_.csv	2020-09-01 오후 4:56	Microsoft Excel ...	7,831KB
교육_2_.csv	2020-09-01 오후 4:53	Microsoft Excel ...	12,623KB

```
7 from bs4 import BeautifulSoup
8 from urllib.request import Request, urlopen
9 import pandas as pd

11 for i in range(1,11):
12     url = '사토 뉴스' % (i)
13
14     headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML}
15
16     request = Request(url, headers = headers)
17     response = urlopen(request)
18     html = response.read()
19     soup = BeautifulSoup(html, 'html.parser')
20     name = soup.select('#main_content > div.list_body.newsflash_body > ul > li > dl > dt')
21     # duration = soup.select('#main_content > div.list_body.newsflash_body > ul> li > dl > dd > spa
22
23     #score = soup.select('body > div.rankingDetailBody > div > table > tr > td > b'})
24
25
26 for n in name :
27     name_l.append(n.text.replace('\n\n', '').replace('\n\r\n\t|\t|\t|\t|\t|\t|\t', '').replace('r
28     score_l.append(n.a.get('href'))
29
30 for s in score_l:
31     url=s
32     request2 = Request(url, headers = headers)
33     response2 = urlopen(request2)
34     html2 = response2.read()
35     soup2 = BeautifulSoup(html2, 'html.parser')
36     content = soup2.select('#articleBodyContents')
37
38
39 for c in content :
40     duration l.append(c.text.replace("\n\n\n\n\n",'').replace("\n\t\n",""))
```

03. 데이터 수집- 뉴스 크롤링

CSV로 저장된 뉴스데이터

교육 - Excel			
송현우			
어떤 작업을 원하시나요?			
파일	홈	삽입	페이지 레이아웃
수식	데이터	검토	보기
도움말			
붙여넣기	맑은 고딕 11	가	가
클립보드	가	가	가
글꼴	맞춤	표시 형식	스타일
스타일	스타일	스타일	스타일
셀	셀	셀	셀
E8			
	A	B	C
1	제목	내용	분류
2	유은혜 "수도권 학교 3단계 미리 준비"...조희연 "원격 전환하	수도권 교육감·기초단체장과 방역점검회의...3단계	교육
3	연세대, 코로나19 확진 대학생 접촉자 17명 전원 음성	연세대 공학원, 코로나19 확진자 발생으로 폐쇄(서	교육
4	수도권 교육감들 "3단계 전 학교 문닫아야"...교육부 "선제조	유은혜 "거리두기 3단계 전제로 등교문제 협의"오	교육
5	교육수장들 "원격수업? 전면 전환 선제적 시행 필요"	유은혜 부총리, 시도교육감 수도권 학교방역 점검	교육
6	브랜섬홀 아시아, 전 세계 우수 대학 대거 참가 비대면 '세계	제주영어교육도시 내 여성 IB(International	교육
7	경북대병원장에 김용림 신장내과 교수	(대구=연합뉴스) 김선형 기자 = 경북대학교병원은	교육
8	경북교육청, 유·초·중·고 학생 밀집도 최소화한다	이달 24일~26일 시작해 9월 11일까지경사북도교	교육
9	(주)필터레인, 한신대 창업보육센터 졸업기념 대학발전기금 코	(주)필터레인 최승우 대표(좌)가 연규홍 총장(우)에	교육
10	'제자 성추행' 혐의 서울대 음대 교수, 5년 만에 불구속 기소	서울대 음대 교수가 차 안에서 제자를 성추행한 혐	교육
11	(주)코리아메디컬, 우즈벡 3개 의대 한국 캠퍼스 설립 계약	한국 사회에서 '의사'는 누구에게나 선망받는 상위	교육

A blurred office scene with a person working at a desk. The desk has multiple monitors, a desk lamp, and various office supplies. The text "데이터 전처리" is overlaid in the center.

데이터 전처리

04. 데이터 전처리 - 잡음제거

제목	내용
유은혜 "수도권 교육감·기초단체장과 방역점검회의...3단계서는 등교 중단, 원격수업 또는 휴업서울교총 "등교수업, 브랜섬홀 전면 온라인 수업으로 전환해야"인사말 하는 유은혜 교육부 장관(서울=연합뉴스) 정하종 기자 = 24일 오후 경북대병원 서울 서대문구 서울시교육청에서 열린 수도권 학교 방역 점검 회의에서 유은혜 사회부총리 겸 교육부 장관 (주필터레이 인사말을 하고 있다. chc@yna.co.kr (서울·세종=연합뉴스) 고유선 김수현 기자 = 유은혜 부총리 겸 교육부 장관은 24일 "수도권 지역의 사회적 거리두기 3단계가 언제라도 나올 수 있는 상황임을 전제하고 교육청과 함께 미리미리 필요한 준비를 하겠다는 말씀도 드린다"고 밝혔다. 유 부총리는 이날 서울시교육청에	<pre>// flash 오류를 우회하기 위한 함수 추가 function _flash_removeCallback() {}</pre>

■ 기사 내용과 관련 없는 내용 제거

- 공백, // flash 오류 우회 함수 추가 → `function _flash_removeCallback() {}`
- 특수문자 , 이메일 크롤링

04. 데이터 전처리 - 형태소 분석

- 기사 분류에 관련 없는 동사, 형태소를 제거하고 명사만을 추출하여 정확한 기사 분류에 도움을 줌

다람쥐 **현** 쳇바퀴에 타고 있다.

다람쥐

쳇바퀴

04. 데이터 전처리 - 형태소 분석

기사내용

온라인으로 진행, 1500명 이상 참여 예정(지디넷코리아=방은주 기자)한국심리학회(회장 조현섭 총신대학교 교수)는 2020 제74차 한국심리학회 연차학술대회를 20~22일 사흘간 온라인으로 개최한다고 19일 밝혔다. 이번 학술대회는 1개 프리세션과 대외 심포지엄, 6개 세션 특별 심포지엄, 11개 세션 분과 심포지엄, 6개 세션 분과 워크숍, 청소년 심리학교실, 한중일 심포지엄 등 총 111편의 포스터 발표가 온라인으로 열린다. 1500명 이상이 참가할 예정이다. 주최 측은 참가 인원수 제한에도 200명이 넘는 중·고등학생이 청소년심리학교실에 참가 신청을 했다고 밝혔다. 1일차(20일)에 개최할 프리세션은 '국민을 위한 심리서비스 법제화의 현실적 필요성과 실현 방안'이라는 주제로 국내 심리서비스 운영의 문제점을 진단하고 개선 방향을 제안한다. 이후 대외심포지엄에서 '위기의 지구: 재난, 심리학에서 길을 찾다'는 주제로 최근 발생한 코로나19 대유행을 중심으로 재난이 우리의 삶과 심리에 미치는 부정적 영향을 이해하고, 이에 대한 대응 방안을 모색하는 자리를 마련했다. 특히, 심리학을 기반으로 의학, 보건학, 의료사회학 등 다양한 인접 학문 관점에서 재난의 피해를 해결하는데 도움이 되는 연구 성과를 논의할 예정이다. 또한 미국심리학회를 대표해 산드라 L. शुल्만(Sandra L. Shullman) 회장이 발표자로 참여한다. 2일차(21일)에는 심포지엄과 워크숍을 통해 우리 사회의 다양한 문제를 다룬다. 한국형 심리방역체계 구축, 중독, 자살, 심리적 외상, 디지털 성범죄, 재범 방지 등과 같은 우리 사회가 직면한 현안들을 심리학적으로 진단하고 그 해결책을 모색한다. 특히 정서장애의 근거 기반치료, 게슈탈트 치료, 코칭 프로그램 등과 같이 우리가 현대사회를 살아

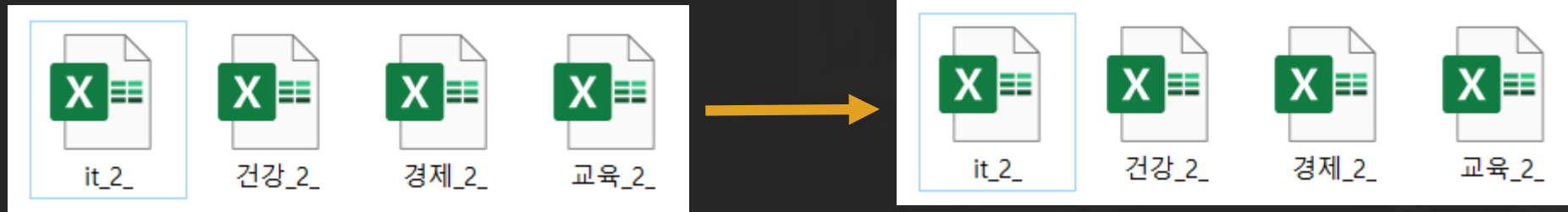
명사 추출

['온라인', '진행', '1500명', '이상', '참여', '예정(지디넷코리아=방은주', '기자)한국심리학회(회장', '조현섭', '총신대학교', '교수', '2020', '제74차', '한국심리학회', '연차학술대회', '20', '22일', '사흘간', '개회', '19일', '이번', '학술대회', '1개', '프리세션', '대외', '심포지엄', '6개', '세션', '특별', '11개', '분', '워크숍', '청소년', '심리학교', '한·중·', '등', '111편', '포스터', '발표가', '참가', '예정', '주최', '측은', '인원수', '제한', '200명', '중·고등학생', '청소년심리학교실', '신청', '1일차(20일)', '국민', '심리서비스', '법제화', '현실적', '필요성', '실현', '방안', '주제', '국내', '운영', '문제점', '진단', '개선', '방향', '제안', '이후', '대외심포지엄', '위기', '지구', '재난', '심리학', '찾다', '발생', '코로나19', '대유행', '중심', '우리', '삶', '심리', '부정적', '영향', '이해', '이', '대응', '방안', '모색', '자리', '마련', '특히', '기반', '의학', '보건학', '의료사회학', '다양한', '인접', '학문', '관점', '피해', '해결', '도움', '연구', '성과', '논의', '미국심리학회', '대표', '산드라', '술먼(Sandra)', '회장', '발표자', '2일차(21일)', '사회', '문제', '한국형', '심리방역체계', '구축', '중독', '자살', '심리적', '외상', '디지털', '성범죄', '재범', '방지', '직면', '현안들', '심리학적', '해결책', '정서장애', '근거', '기반치료', '게슈탈트', '치료', '코칭', '프로그램', '현대사회', '다양', '맥락', '문제들', '실용적', '프로그램들', '다수', '준비', '마지막', '3일차(22일)', '청소년들', '응용', '수',



모델 학습

05. 모델학습 – 데이터 쌓기



형태소 분석이 끝난 파일

카테고리 정수화 : 0 1 2 3

```
1 w_list = []
2 for item in data_0.내용:
3     item = item[1:-1]
4     temp = item.split(',')
5     w_list.append(temp)
6
7 w_list_result = []
8 for n in w_list:
9     news = []
10    for x in n:
11        x = x.strip()
12        if(len(x) > 1):
13            news.append(x)
14    w_list_result.append(news)
15    y_category.append(0)
16
17 #IT Append

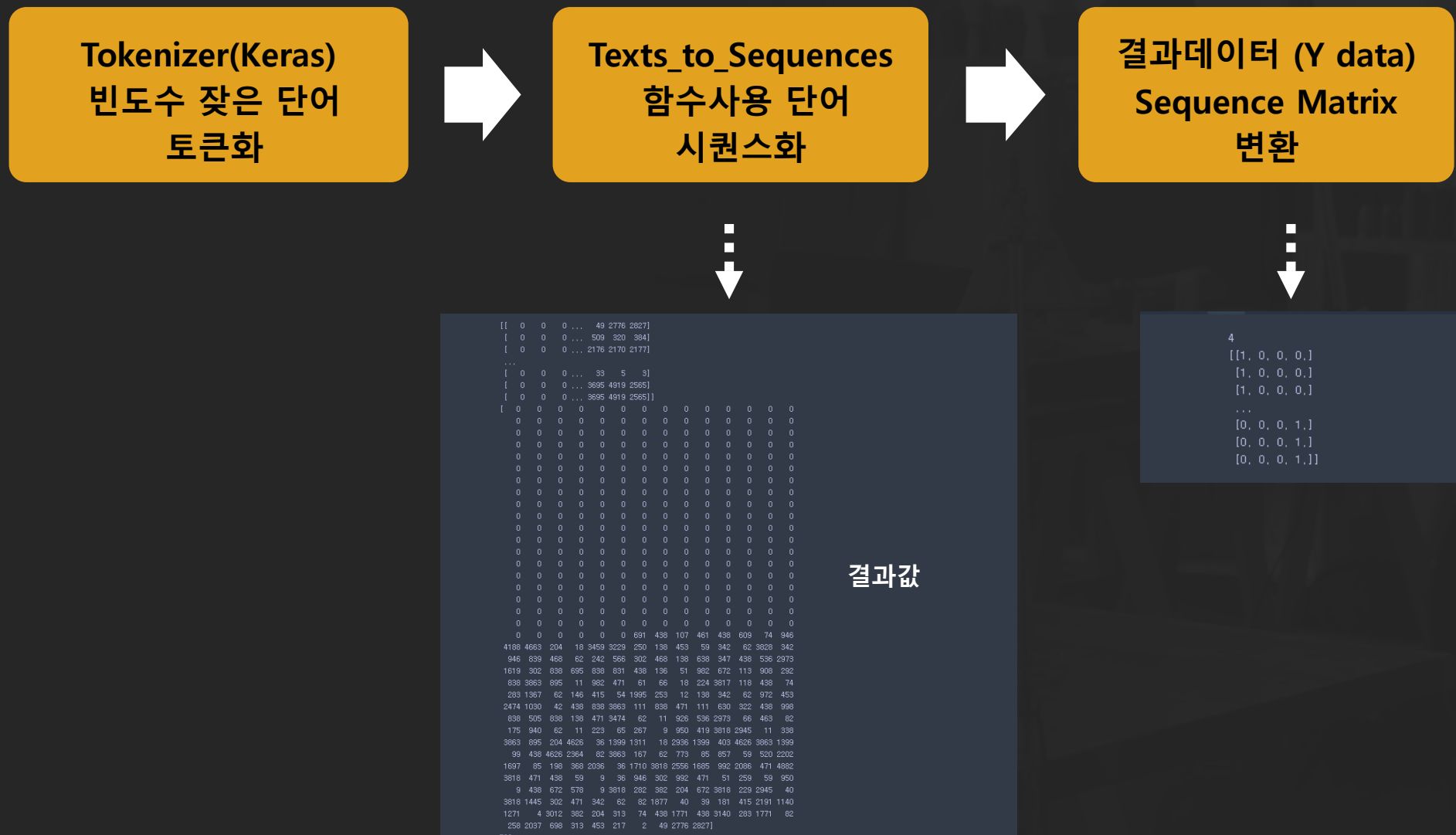
1 w_list = []
2 for item in data_1.내용:
3     item = item[1:-1]
4     temp = item.split(',')
5     w_list.append(temp)
6
7 for n in w_list:
8     news = []
9     for x in n:
10        x = x.strip()
11        if(len(x) > 1):
12            news.append(x)
13    w_list_result.append(news)
14    y_category.append(1)
15
16 #건강 Append

1 w_list = []
2 for item in data_2.내용:
3     item = item[1:-1]
4     temp = item.split(',')
5     w_list.append(temp)
6
7 for n in w_list:
8     news = []
9     for x in n:
10        x = x.strip()
11        if(len(x) > 1):
12            news.append(x)
13    w_list_result.append(news)
14    y_category.append(2)
15
16 #경제 Append

1 w_list = []
2 for item in data_3.내용:
3     item = item[1:-1]
4     temp = item.split(',')
5     w_list.append(temp)
6
7 for n in w_list:
8     news = []
9     for x in n:
10        x = x.strip()
11        if(len(x) > 1):
12            news.append(x)
13    w_list_result.append(news)
14    y_category.append(3)
15
16 #교육 Append
```

입력 데이터 저장 (X data)
결과 데이터 저장 (Y data) 저장

05. 모델학습 - 토큰화, 시퀀스화



05. 모델학습 - 모델 생성



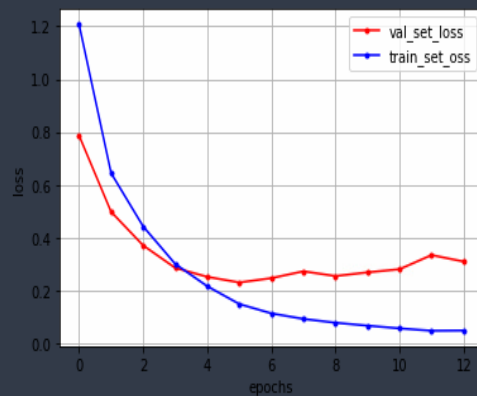
05. 모델학습 - 학습결과

학습 결과
Accuracy: 93%

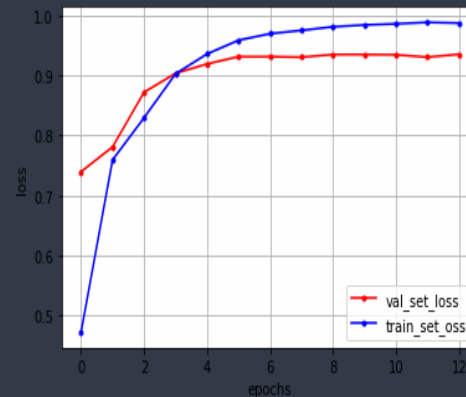
```
In [40]: | print("정확도 : %.4f" % (model.evaluate(X_test, y_test)[1]))
```

223/223 [=====] - 15s 69ms/step - loss: 0.2903 - accuracy: 0.9368
정확도 : 0.9368

Loss Graph



Accuracy Graph



A blurred office scene with a person working at a desk with multiple computer monitors. The text "서비스 개발" is overlaid in the center.

서비스 개발

06. 서비스 개발 – 초기 디자인



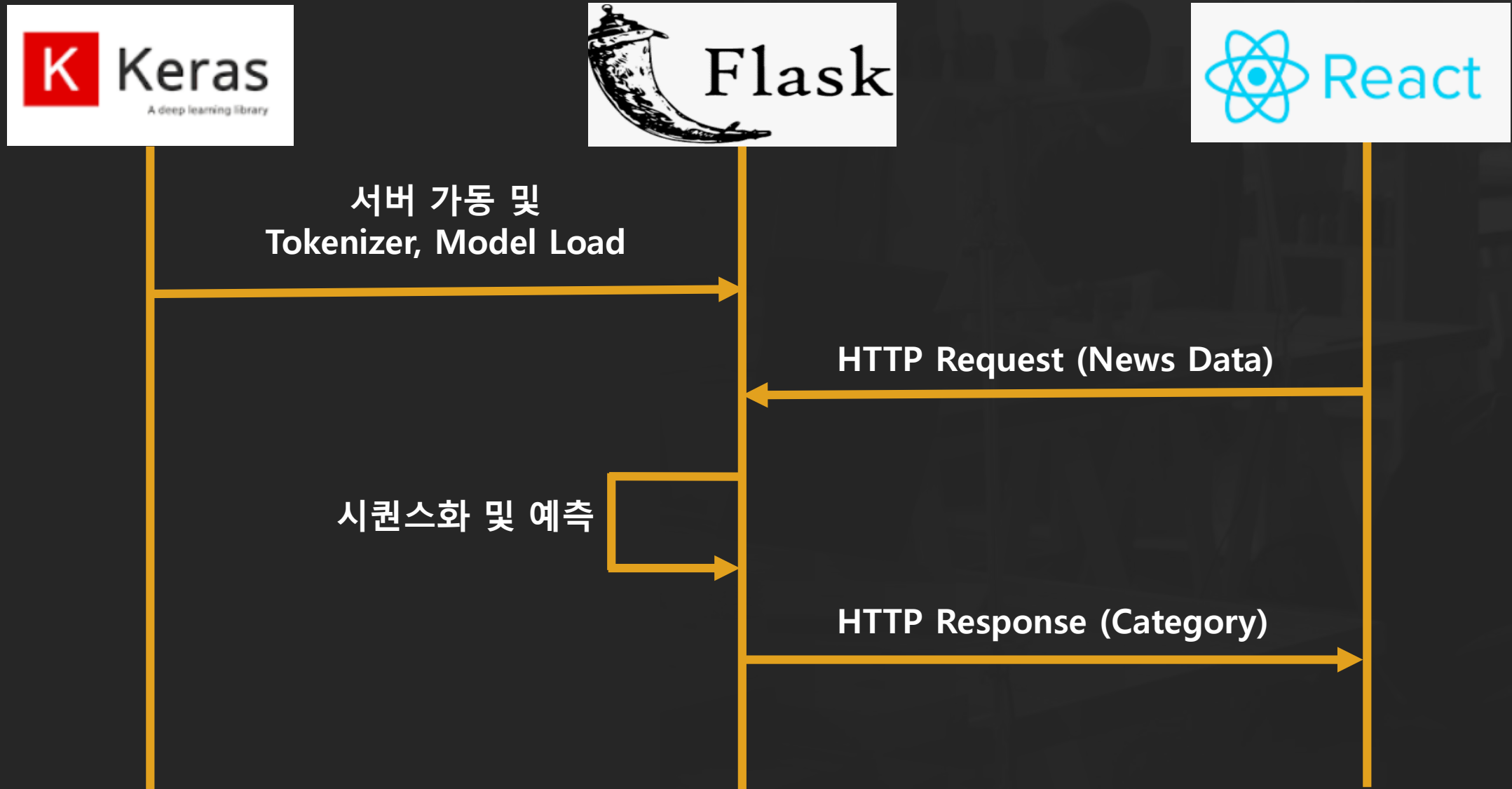
파일선택

기사 리스트

카테고리

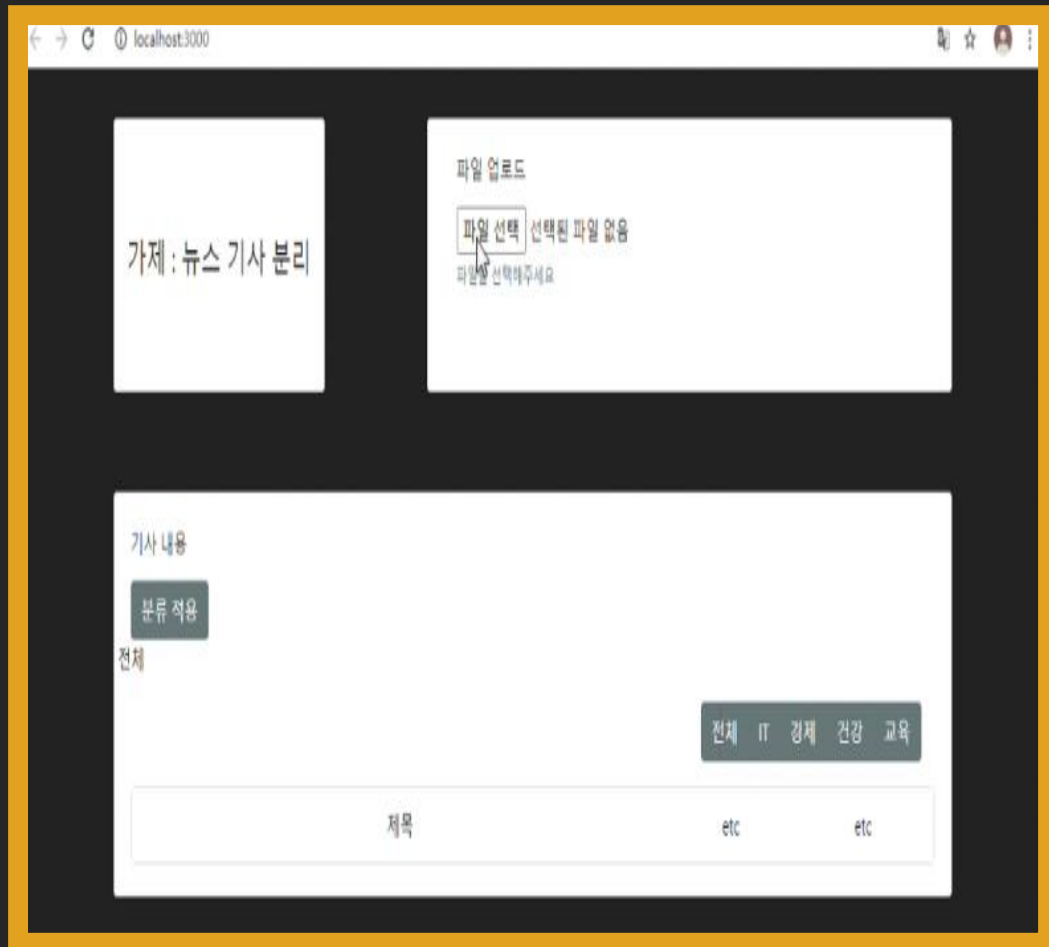
결과 보기

06. 서비스 개발 – Back-End 서버

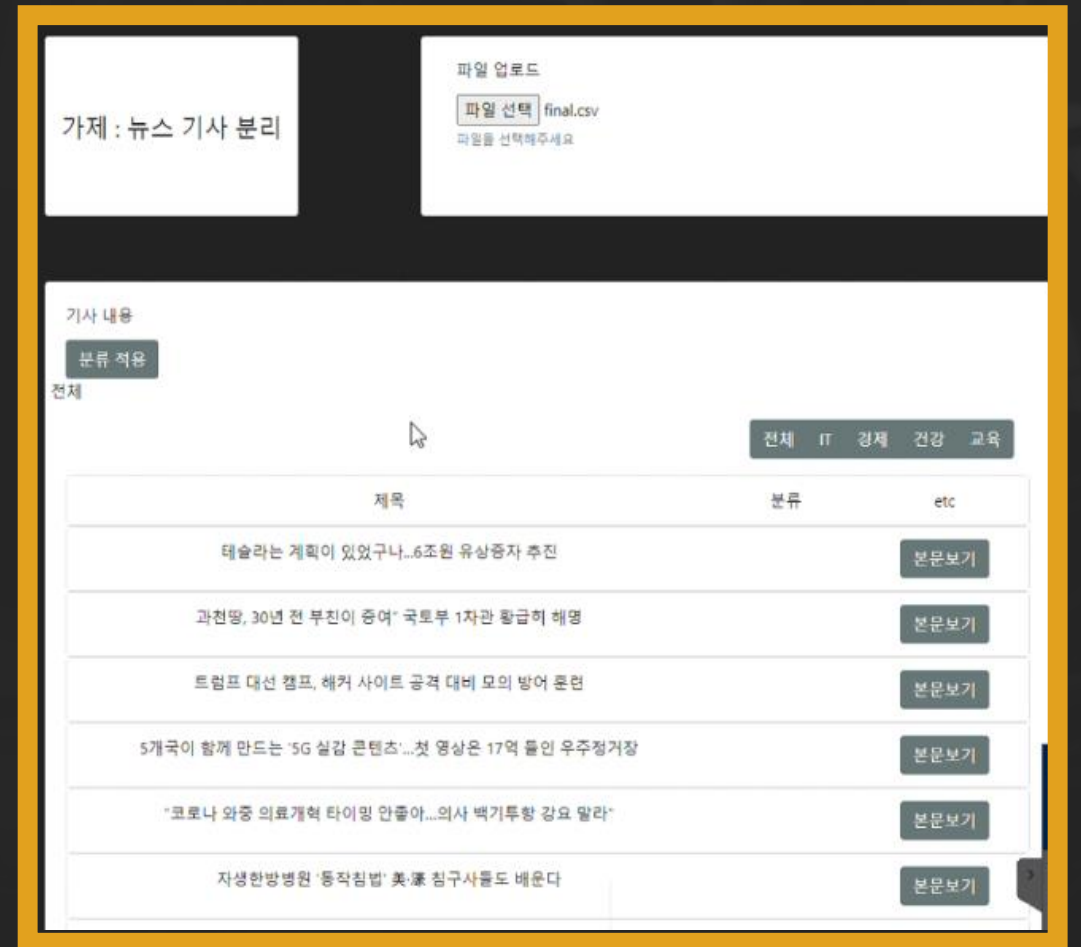


06. 서비스 개발 - 서비스 사용

(1) 기사 파일 선택



(2) 기사 로드 완료



06. 서비스 개발 – 서비스 사용

(3) 분류 적용 버튼 클릭

가제 : 뉴스 기사 분리

파일 업로드 | titbit.csv
파일을 선택해주세요

기사 내용

전체

분류 적용

전체 IT 경제 건강 교육

제목	분류	etc
테슬라는 계획이 있었구나...6조원 유상증자 추진		본문보기
과천땅, 30년 전 부친이 증여' 국토부 1차관 황급히 해명		본문보기
트럼프 대선 캠프, 해커 사이트 공격 대비 모의 방어 훈련		본문보기
5개국이 함께 만드는 '5G 실감 콘텐츠'...첫 영상은 17억 들인 우주정거장		본문보기

06. 서비스 개발 - 서비스 사용

(4) 분류 완료 및 분야별 모아 보기

분류 적용				
전체	IT	경제	건강	교육
제목	분류	etc		
태슬라는 계획이 있었구나..6조원 이상증자 추진	IT	본문보기		
과천영, 30년 전 부친이 중여"국토부 1차관 황급히 해명	경제	본문보기		
트럼프 대선 캠프, 해커 사이트 공격 대비 모의 방어 훈련	경제	본문보기		
5개국이 함께 만드는 '5G 실감 콘텐츠'..첫 영상은 17억 돌인 우주정거장	IT	본문보기		
'코로나 와중 의료개혁 타이밍 안좋아...의사 벼기듯할 강요 말라'	IT	본문보기		
자생전병병원 동적침방 美 藥 침구사들도 배운다	건강	본문보기		
혈변, 복통... '아래의 병' 증상성장질환 의심	건강	본문보기		
폐암 등 3개질환 동시발목...AI SW 나온다	건강	본문보기		

파일 업로드

파일 선택

final.csv

파일을 선택해주세요

가제: 뉴스 기사 분리

기사 내용

분류 적용

IT

테슬라는 계획이 있었구나..6조원 유상증자 추진

IT

본문보기

올들어 주가가 500% 급등한 미국 전기자동차 업체 테슬라가 유상증자(capital raise)를 통해 50억달러(약 5조9200억원)를 조달하기로 했다. 골드만삭스,뱅크오브아메리카, 버클레이스, 피티코, 모건스탠리 등 다수의 금융회사들을 통해 수시로 신주를 발행할 계획이다. 1일(현지시간) 미국 언론들은 "일론 머스크 CEO(최고경영자)는 영업을 통해서도 충분히 성장할 수 있기 때문에 신규 자금 조달에는 관심이 없다고 선을 그어 왔지만, 최근 주가 급등으로 유상 증자 유혹에 휩싸인 것 같다"고 보도했다. 통상 유상증자는 재무 구조가 취약한 상장기업이 주가가 상승할 때 활용하는 자금조달 수단이다. 주가 급등기에는 같은 수의 주식을 발행해도 유입되는 자금이 더 많기 때문이다. 가령 주가가 500원일 때 100주를 발행하면 5만원이 들어오지만, 5000원일 때는 50만원이 유입되어 회사 입장에선 더 유리하다. 테슬라는 지난 2월에도 주가 급등기에 20억불을 조달한 바 있다. 당시 테슬라는 대차대조표 강화 목적의 재원으로 활용할 계획이라고 발표한 바 있다. 당시 대규모 유상증자여도 불구하고 테슬라 주가는 계속 상승해 이천달러(주가 2000불 돌파)라고 불릴 정도로 크게 상승했다. 한편 예상 밖 대규모 유상증자 추진 소식이 전해지자, 가장 천 거액에서 7% 상승하며 500달러를 넘어섰던 테슬라 주가는 4%대 하락한 475달러 범위 내 하락 모습이다. 한편, 5대1 역분할 이후 지난 달 31일(현지시간) 첫 거래를 시작한 테슬라는 전일 대비 12.6% 저수어 498.32달러로 마감하면서 시가총액이 4643억달러를 기록했다. 이는 세계 최대 유통업체인 월마트의 시총(3932억달러)과 세계 최대 카드사인 비자(4511억달러)의 시가총액을 모두 넘어서는 것이다.

5개국이 한데 만드는 '5G 실감 콘텐츠'...첫 영상은 17억 원의 우주정거장

IT

본문보기

파일 업로드

파일 선택

final.csv

파일을 선택해주세요

가제 : 뉴스 기사 분리

기사 내용

분류 적용

IT

전체

IT

경제

건강

교육

태슬라는 계획이 있었구나..6조원 유상증자 추진	IT	본문보기
5개국이 함께 만드는 '5G 실감 콘텐츠'..첫 영상은 17억 들인 우주정거장	IT	본문보기
"코로나 외종 의료개혁 타이밍 안 좋아...외사 백기투항 강요 할라"	IT	본문보기
사립대 86% '당분간 비대면수업' ..10%는 2학기 내내	IT	본문보기
전북은행 '디지털 특화' 신용카드 출시...유통브 넷플릭스-엘론 최대 1만원 할인	IT	본문보기

06. 서비스 개발 – 검색활용(집중호우)

전체 IT 경제 건강 교육		
제목	분류	etc
집중호우에 농산물가격 경중...소비자물가 두달째 상승세(상보)	경제	본문보기
코로나로 충남 시·군 학교급식지원센터 운영난 심각	경제	본문보기
태풍 '마이삭' 북상 부산 학교 2일 하교 조정, 3일 원격수업 권장	교육	본문보기
장마·호우에 8월 채소값 28.6%↑...소비자물가 5개월만 최대폭 상승	경제	본문보기
8월 소비자물가 작년비 0.7%↑...집중호우로 '농축수산물' 가격 급등	경제	본문보기
구례군, 수해 주민 건강관리 강화	건강	본문보기
대전시, 집중호우 피해·복구 상황 및 방재대책 추진	경제	본문보기
8월 소비자물가 5개월 만에 최대↑...집중호우에 채소류 급등(종합)	경제	본문보기
수해민 대상 방문건강 관리	경제	본문보기
즉석밥에 디저트까지... 코로나·장마에 식품업체 줄줄이 가격 인상	IT	본문보기

2020-09-02.CSV - Excel

파일 홈 삽입 레이아웃 수식 데이터 검토 보기 도움말 어떤 작업을 원하시나요?

잘라내기 붙여넣기 서식 복사 서식 붙여넣기 글꼴 배경색 채우기 테두리 스타일

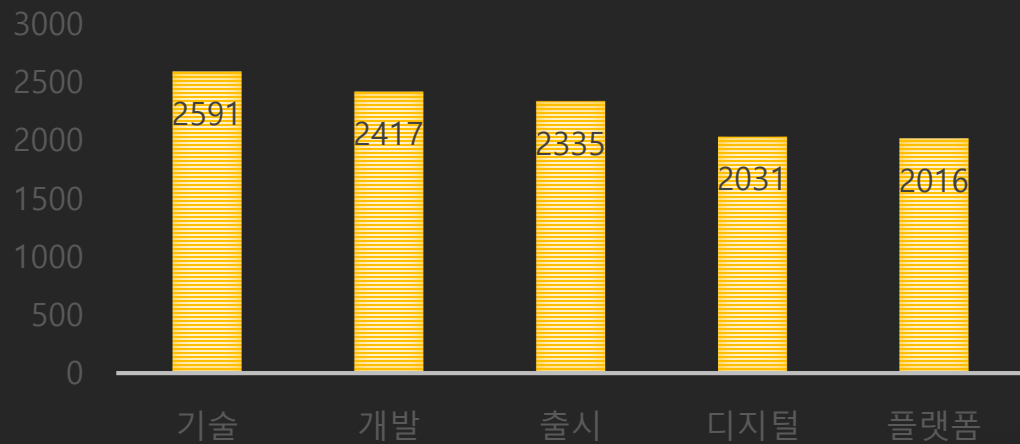
데이터가 손실될 수 있음 이 통합 문서를 심플로 구분된 형식(.csv)으로 저장하면 일부 기능이 손실될 수 있습니다. 기능을 유지하려면 Excel 파일 형식으로 저장하세요.

A	B	분류
1 제목	내용	
2 집중호우에 농산물가격 경중...소비자물가 두달째 상승세(상보)	[세종=이데일리 이명철 기자] 국내 소비자	
3 코로나로 충남 시·군 학교급식지원센터 운영난 심각	코로나19 확산으로 충남지역 학교들의 등교	
4 태풍 '마이삭' 북상 부산 학교 2일 하교 조정, 3일 원격수업 권장	[에듀인뉴스=한치원 기자] 부산시교육청이	
5 장마·호우에 8월 채소값 28.6%↑...소비자물가 5개월만 최대폭 상승	지난달 소비자물가 상승률이 0.7%를 기록했	
6 8월 소비자물가 작년비 0.7%↑...집중호우로 '농축수산물' 가격 급등	지난달 국내 소비자물가가 작년 동월 대비	
7 구례군, 수해 주민 건강관리 강화	구례군이 8월 집중호우로 마을·주택 등이 침	
8 대전시, 집중호우 피해·복구 상황 및 방재대책 추진	대전시는 지난 7월 30일 시간당 79mm의 집	
9 8월 소비자물가 5개월 만에 최대↑...집중호우에 채소류 급등(종합)	[세종=뉴스1] 박영주 기자 = 지난달 소비	
10 수해민 대상 방문건강 관리	구례군이 이번 집중호우로 마을, 주택 등이	
11 즉석밥에 디저트까지... 코로나·장마에 식품업체 줄줄이 가격 인상	오투기 즉석밥 8% 인상, 롯데제과·칠성음료	

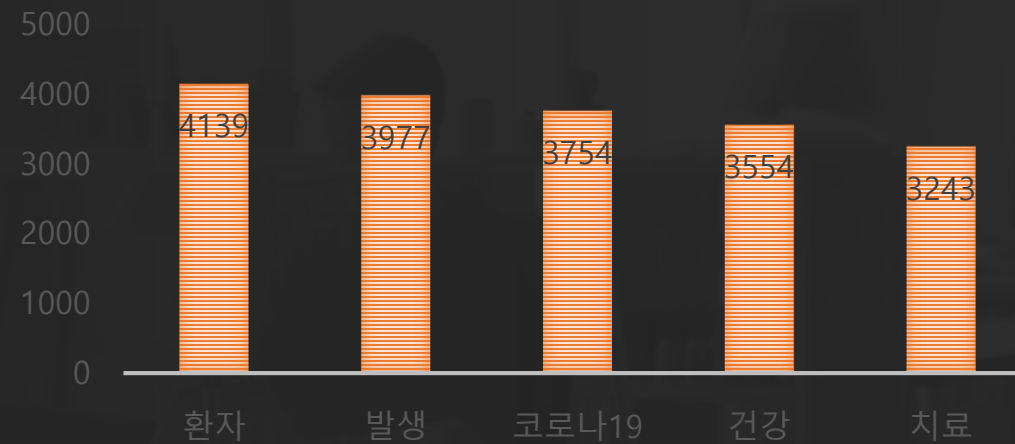
2020-09-02

네이버 검색 키워드 “집중호우”
뉴스 리스트

IT



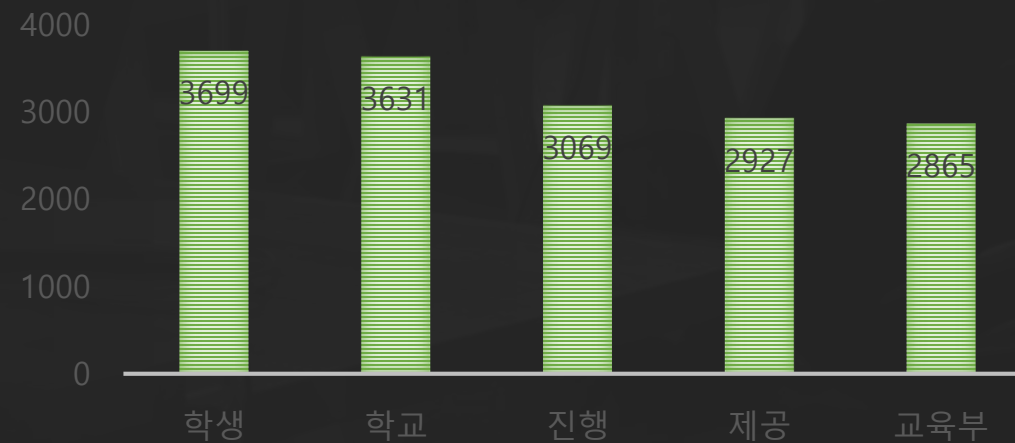
건강



경제



교육





활용 방안 및 기대사항

07. 활용 방안 및 기대 사항

검색엔진

포털 검색엔진에 이용 가능

자동화

- 지식인/블로그에 해시태그 자동화 기능 추가 가능
- 세분화 되고 전문성 있는 정보 획득

연관성

키워드 검색 시 분야별 기사
포커스 와 입장 등을 알 수 있다

A blurred office scene with a person working at a desk in the background and various office supplies on a desk in the foreground. The text "Q&A" is centered in the middle of the image.

Q&A



감사합니다.