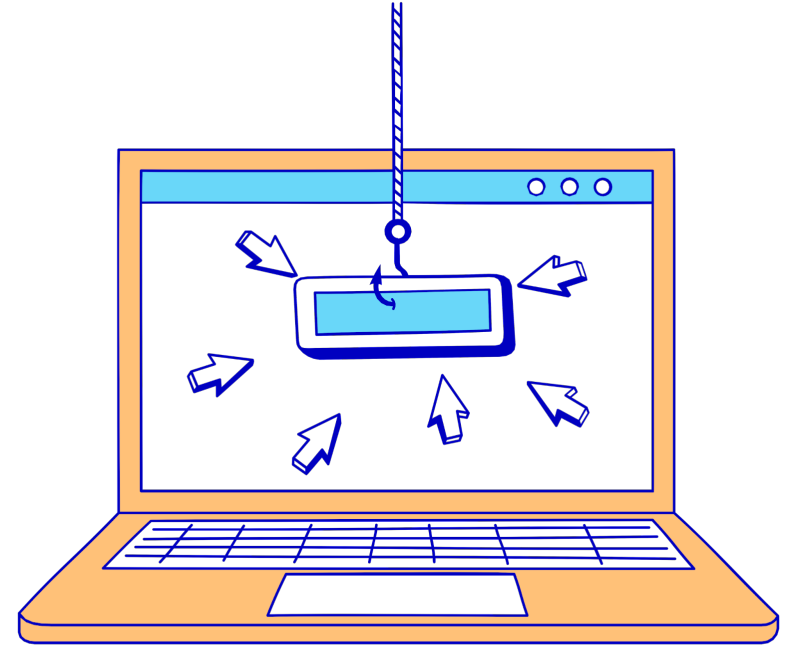


딥러닝 모델을 통한 클릭베이트 기사 감지



1조- 데이터구조

김예담

김지예
우형석

김현민
이한

목차

1

- 팀구성
- 프로젝트 배경

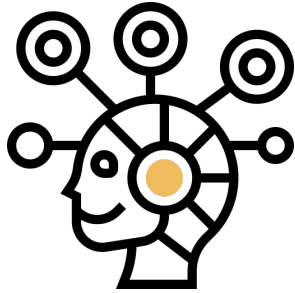
2

- 딥러닝 모델 구축
- 백엔드 & 환경구축
- 크롤링 & 프론트엔드

3

- 결과물 시현
- 소감
- 참고문헌
- Q&A

1. 팀구성



김현민

딥러닝 모델 개발



김예담, 김지예

환경구축
워드클라우드
Crontab
DB기획

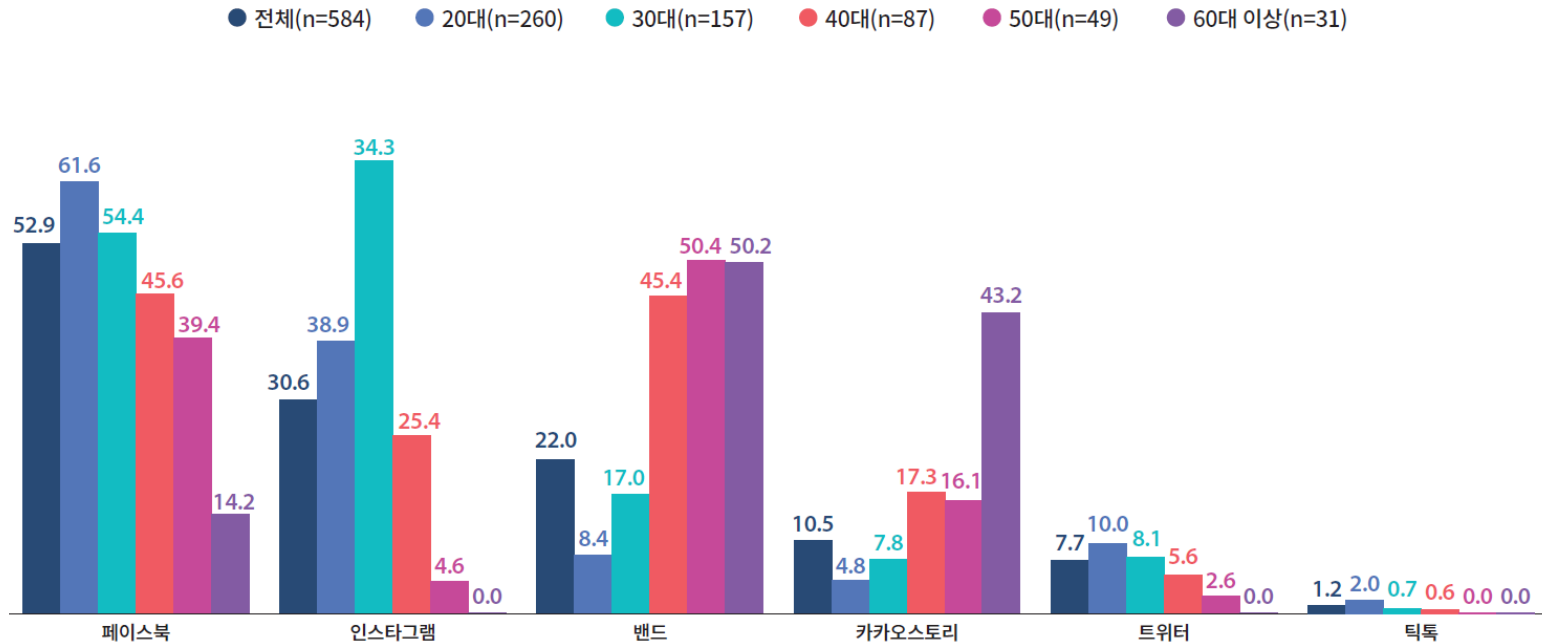


우형석, 이한

크롤링
프론트

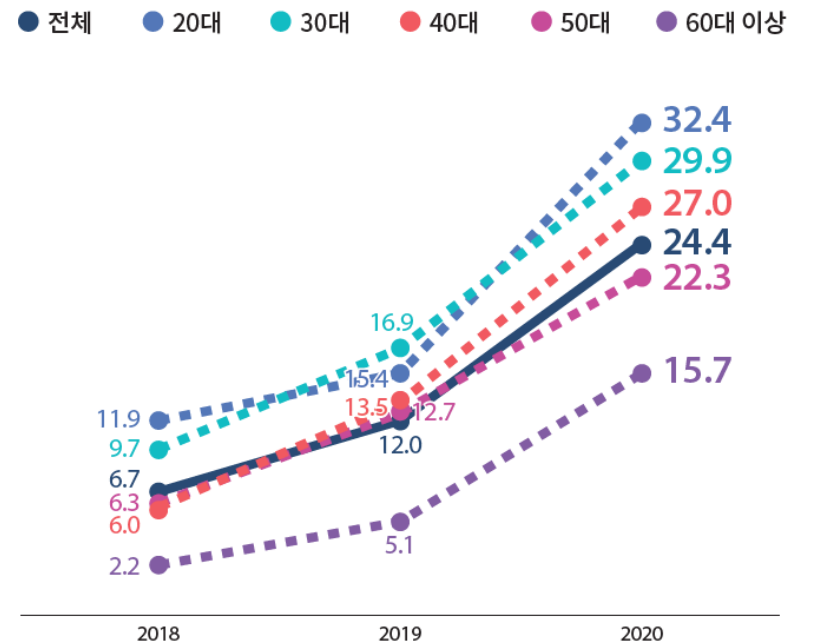
1. 프로젝트 배경

SNS 뉴스 이용자의 SNS별 뉴스 이용률 (단위: %)



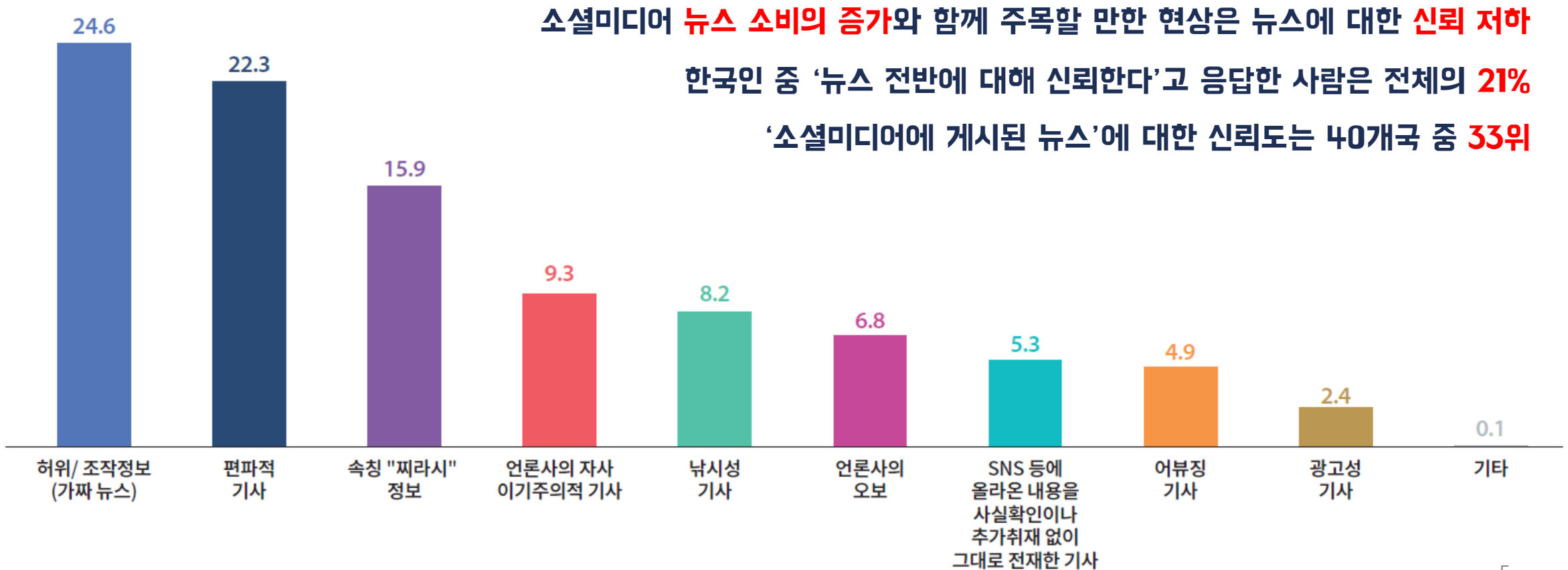
한국언론진흥재단의 <2020년 언론수용자 조사>

연령대별 온라인 동영상 플랫폼을 통한 뉴스 이용률 추이(2018~ 2020년) (단위: %)



- 2020년 유튜브 등 온라인 동영상 플랫폼을 통한 뉴스 이용률은 **24.4%**로 최고치를 기록
- 페이스북을 비롯한 SNS, 메신저, 유튜브를 통한 뉴스 소비가 보편적으로 자리 잡아가는 상태.

한국 언론의 가장 큰 문제점 (단위: %)



클릭베이트 기사 탐지용 AI 모델 개발

클릭베이트 기사 데이터 구축



Annotators

제목 기반 자극성 판단

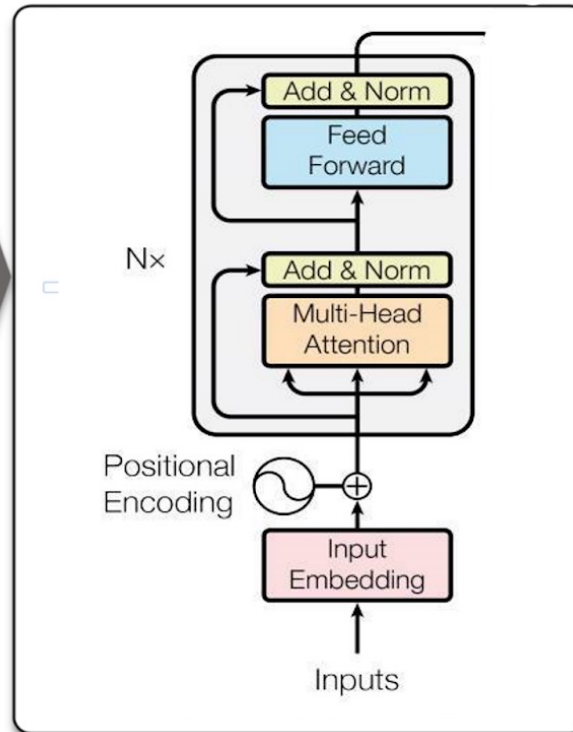


제목과 본문의 불일치 기사



클릭베이트 기사 탐지 한국어 모델

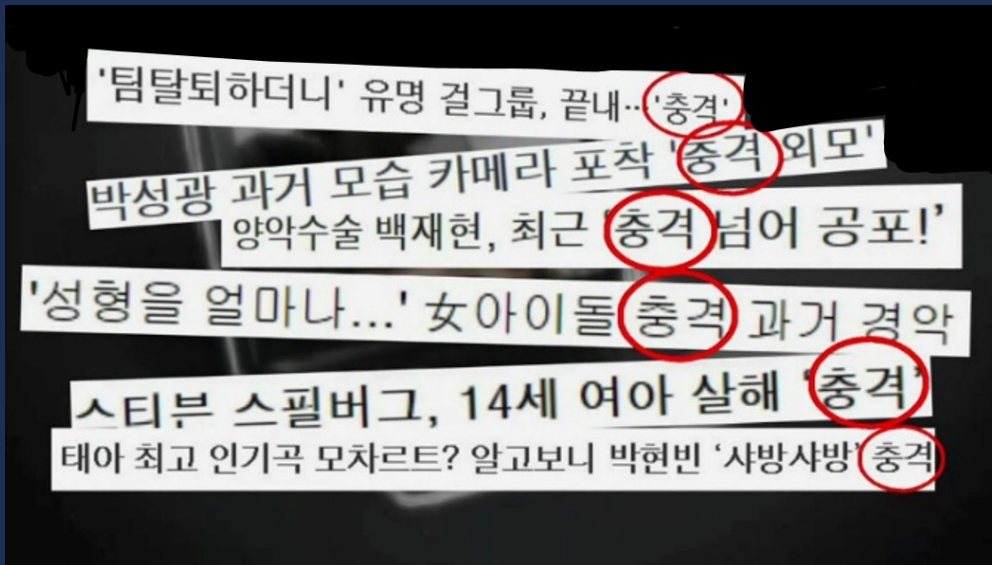
- 두 가지 목적에 부합하는 모델 개발
- 낚시성 기사 탐지 모델의 기준 성능 제시



- 독자들에게 뉴스 제목과 본문의 일관성, 제목의 낚시성 등을 예측하고 점수를 제시하는 **클릭베이트 기사 탐지 시스템** 및 웹사이트 구현

- 이러한 시스템은 올바른 뉴스 기사를 선택할 수 있도록 하여 바람직한 뉴스 생태계를 만드는 데 기여

클릭베이트(ClickBait)란?



- 클릭베이트는 '클릭(Click)'과 '미끼(Bait)'의 합성어
- 국내 클릭베이트는 소위 말하는 '낚시 기사'나 '쓰레기 기사'의 형태로 널리 알려져 있다.
- '헉', '충격', '경악', '역대급' 등의 지나치게 감정적인 표현이나 '숨 막히는 뒤탄' 같은 선정적인 표현을 동반

DATA

모델 학습용 기사 데이터

title	content	label
전주시, 부동산 불법행위 신고 1000만원 포상	전북 전주시는 단속 사각지대에서 음성적으로 이뤄지는 부동산 거래 불법행위 근절을 위...	1
한기영 서울시의원, '서울청년 선거 대책본부' 구성	한기영 서울시의원이 오는 4월 7일 치러지는 보궐선거에 청년들의 투표 독려를 위해 ...	1
대구시·시장 평가 급상승	지난달 대구시의 주민생활만족도와 시장에 대한 평가가 급상승 한 것으로 나타났다 대 구...	1
北 영변 핵시설서 의문의 연기... "재가동 징후"	북한이 영변 핵시설 단지의 일부 시설을 재가동했다는 주장이 나왔다 단지 내 방사화학...	1
"몸 상태 최악..." 19살에 암 투병 고백한 '스걸파' 클루씨 리더 이채린 근황	지난 1월 종영한 스트릿댄스 걸스 파이터에 댄스 크루 클루씨 리더로 출연한 이채...	0

데이터 출처

- a. AI Hub 뉴스 기사 기계독해 데이터 (약 4만개)
- b. 인터넷 언론사 크롤링 데이터 (약 2만개)

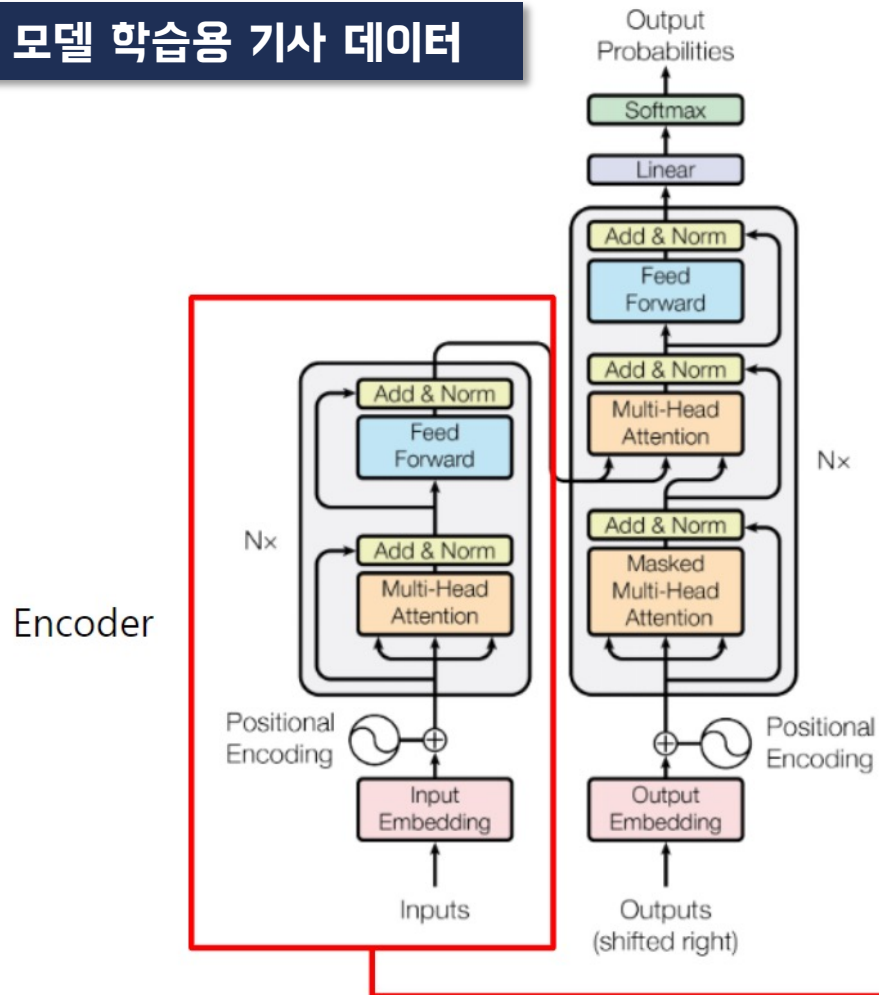
데이터 형태: 제목(title), 내용(content), 라벨

전처리: 중복제거, 특수문자 제거 등

최종 학습 데이터 약 **5만개**

Transformer and BERT

모델 학습용 기사 데이터



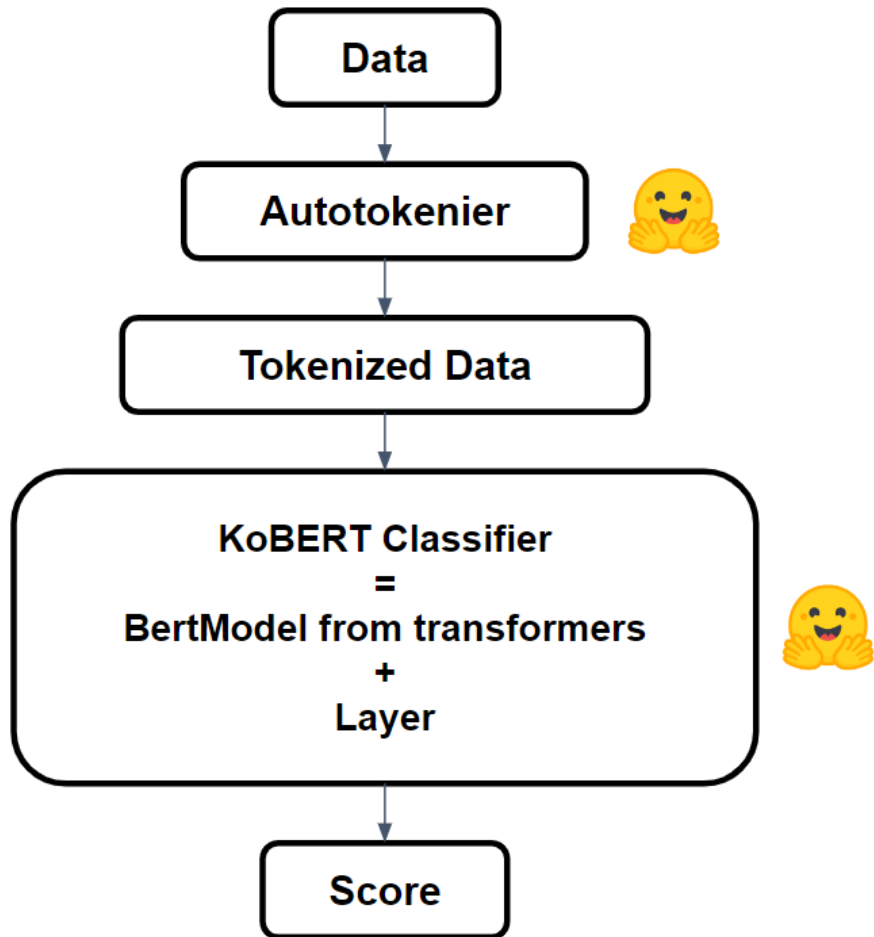
Transformer

- Long-term Dependency 문제 해결

BERT

- Transformer의 Encoder 구조 사용
- 사전 학습 모델
- 양방향 모델을 적용하여
- 문장의 앞과 뒤의 문맥을 고려

KoBERT Classifier



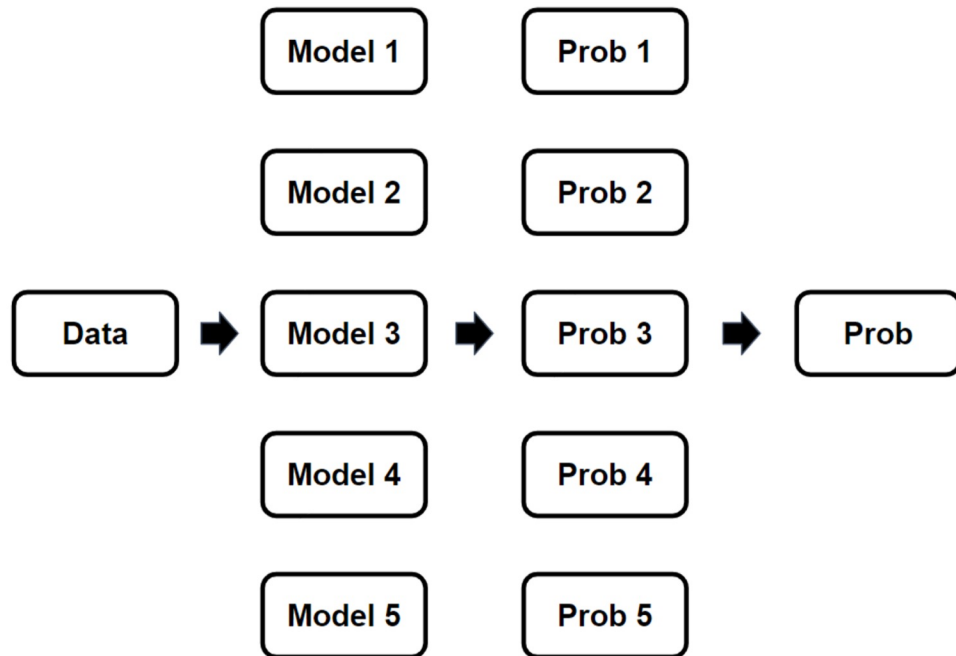
- 데이터 -> 사전학습 Tokenizer -> Bert Input (Tensor 형태)
- 사전학습 KoBERT Model
- BertModel의 Output이 786이기 때문에 추가 Layer 사용, Binary Classifier 정의

Training

- Stratified K-Fold : 5 Fold
- 10 Epoch
- Early stopping
- Optimizer : AdamW
- Scheduler : Linear
- Loss : Label Smoothing CrossEntropy Loss
- Score : F1 Score

Predict Score

모델 학습용 기사 데이터

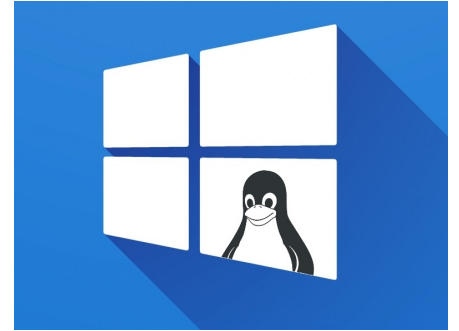


점수 산출 방식

- 정상 기사를 **1**, 비정상 기사를 **0**
- 학습된 모델을 사용한 예측 결과가 **0일 확률** = 해당 기사가 클릭베이트 기사일 확률

점수는 5개 모델의 **예측값 평균** 사용
(Soft Voting)

환경구축



작업 스케줄링(Crontab)

```
# CRONTAB_DJANGO_SETTINGS_MODULE= 'finalProject.settings'

CRONJOBS = [
    ('0 9 * * *', 'article.cron.cronCrawling', '>>> '+os.path.join(BASE_DIR, 'config/log/cron2.log')),
]
```

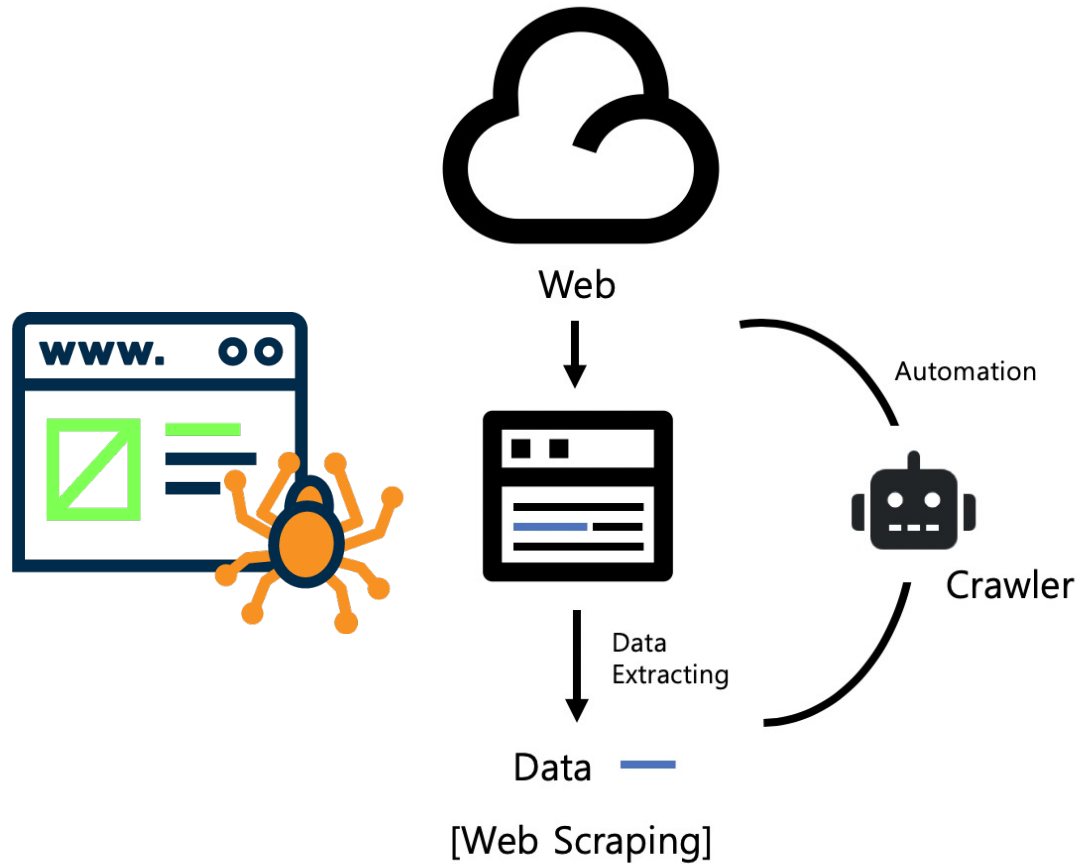
```
waterclean101@DESKTOP-BG91V60: /mnt/c/venvs/final/wslFinal/finalProject
(vfinal) waterclean101@DESKTOP-BG91V60:/mnt/c/venvs/final/wslFinal/finalProject$ python manage.py crontab add
current device : cpu
Seed set as 42
adding cronjob: (f35bdac6428cb94c9a5bfbacb7bfb1c5) -> ('0 9 * * *', 'article.cron.cronCrawling', '>>> /mnt/c/venvs/final/wslFinal/finalProject/config/log/cron2.log')
(vfinal) waterclean101@DESKTOP-BG91V60:/mnt/c/venvs/final/wslFinal/finalProject$ sudo service cron start
[sudo] password for waterclean101:
* Starting periodic command scheduler cron [ OK ]
(vfinal) waterclean101@DESKTOP-BG91V60:/mnt/c/venvs/final/wslFinal/finalProject$ crontab -i
0 9 * * * /mnt/c/venvs/final/wslFinal/bin/python /mnt/c/venvs/final/wslFinal/finalProject/manage.py crontab run f35bdac6428cb94c9a5bfbacb7bfb1c5 >> /mnt/c/venvs/final/wslFinal/finalProject/config/log/cron2.log # django-cronjobs for config
(vfinal) waterclean101@DESKTOP-BG91V60:/mnt/c/venvs/final/wslFinal/finalProject$
```

워드클라우드란?



- 자료의 빈도 시각화
- 빅데이터(big data)를 분석할 때 데이터의 특징을 도출하기 위해 활용
- KoNLPy(한국어 형태소 파서 라이브러리)사용

Web Crawling



목표

- 딥러닝 학습과 웹페이지의 결과물 출력을 위한 뉴스 기사 내용과 제목 추출
- 백엔드, 프론트엔드, 딥러닝에 필요한 모든 요소 수집

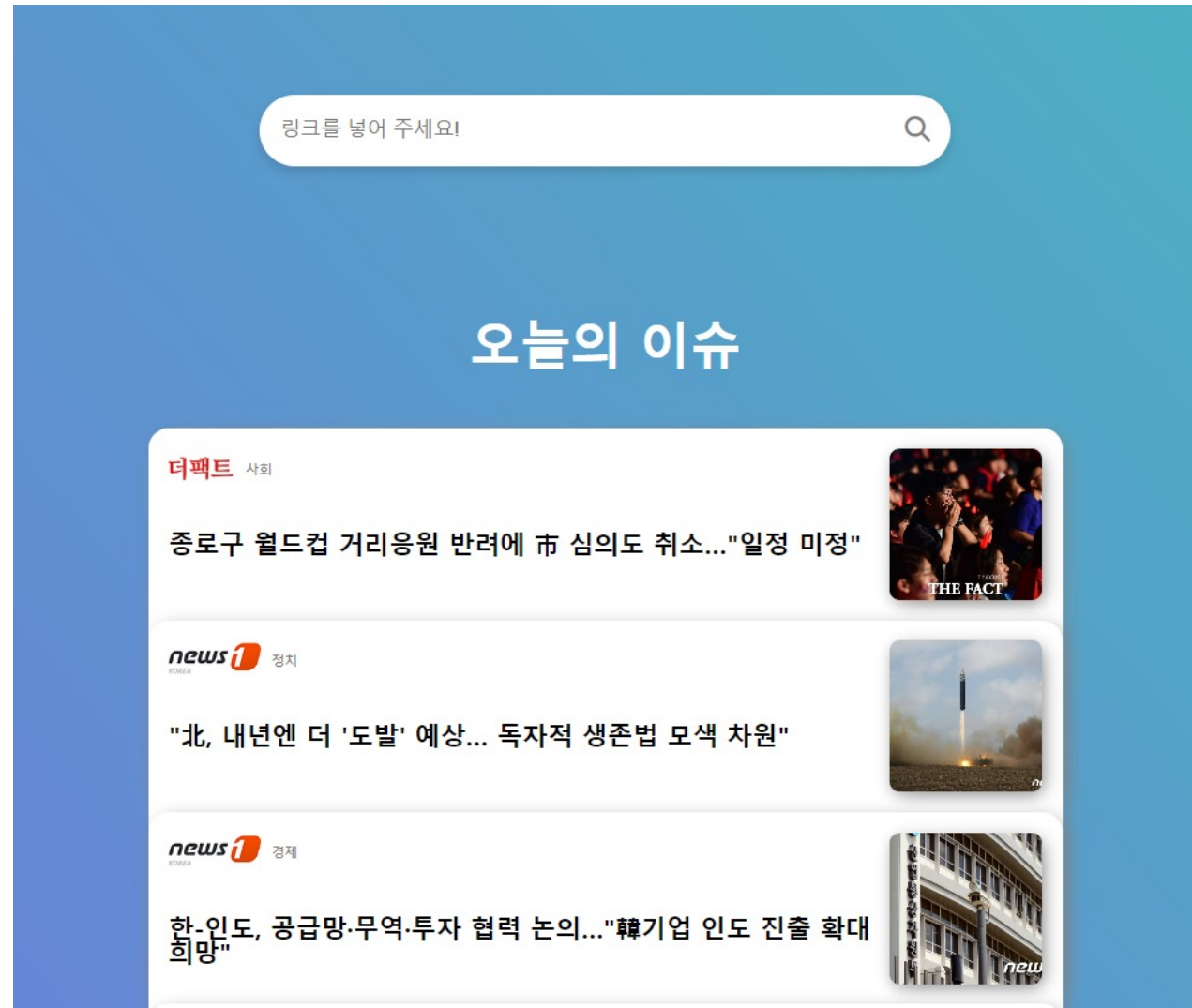
Selenium & BeautifulSoup



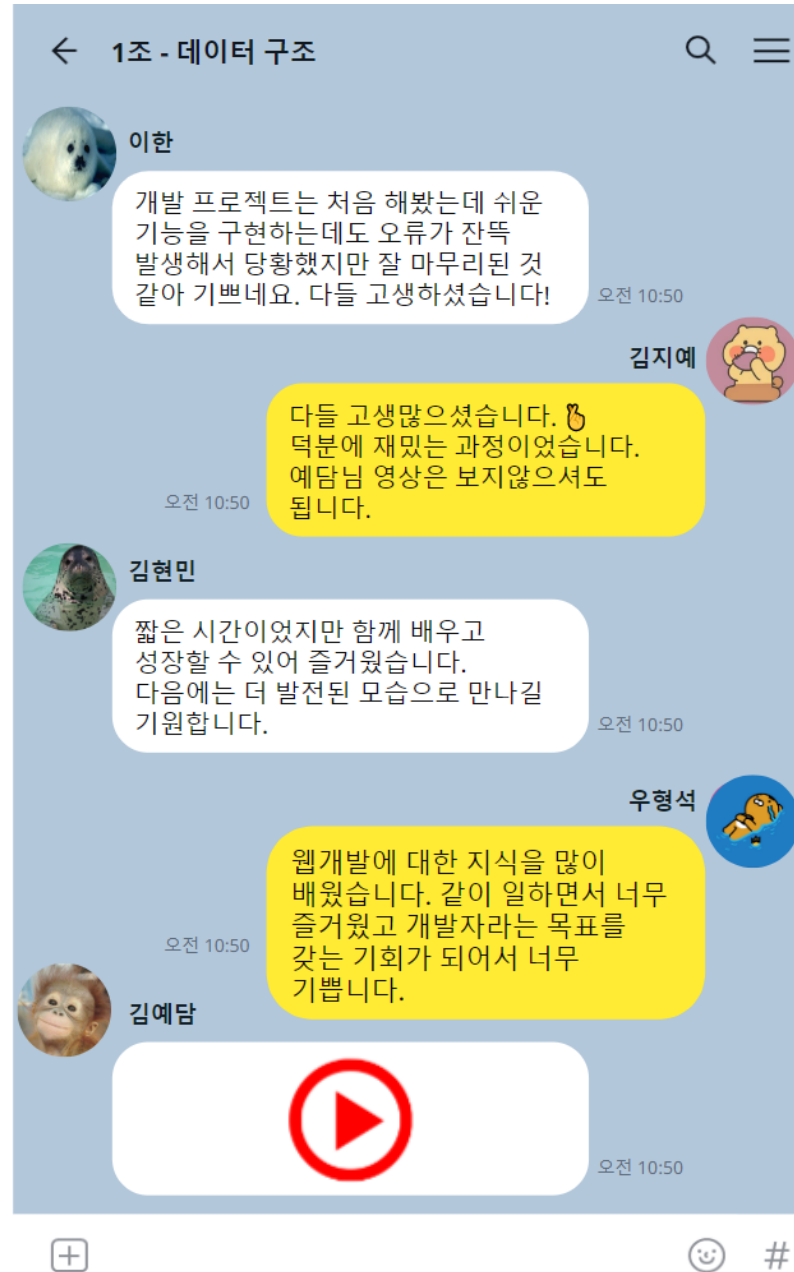
- 자동화된 웹 테스트 프레임워크
- 웹과 동작 기능 자동 실행
- 작동시간 김



- HTML parser
- HTML 요소 추출
- 작동시간 빠름



3. 소감



공식문서

- Transformers Documentation (<https://huggingface.co/docs/transformers/main/en/index>)
- Google BERT - Pre Training and Fine Tuning for NLP Tasks (<https://ranko-music.medium.com/googles-bert-nlp-5b2bb1236d78>)
- Django Documentation

논문

- 2020언론수용자조사, 한국언론진흥재단, 2020.12.15
- 박아란, 이소은, “디지털 뉴스 리포트 : 한국 2020”, 한국언론진흥재단, 2020
- 김태균, 박건우, 차미영, “뉴스 기사의 일관성 탐지를 위한 딥러닝 시스템”, 한국컴퓨터종합학술대회 논문집, 2018, 2201-2203
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin, “Attention Is All You Need”, 2017

Q & A