

인터넷 기사의 댓글 속 혐오 양상에 대하여

20205178 박 현

인공지능 알고리즘을 이용한 사회과학연구 – 기말 프로젝트

최근 온·오프라인 상에서의 혐오가 사회적 문제로 떠오르고 있다. 이에 수 년간의 인터넷 기사 댓글을 수집하여 Smilegate AI에서 공개한 Korean UnSmile Dataset에서 제공하는 인공지능 모델을 이용하여 각 댓글에 해당하는 혐오 분류를 찾아 이들에 대한 분포 및 양상을 분석하고자 한다.

future: Unified Parallel and Distributed Processing in R for Everyone

future.apply: Apply Function to Elements in Parallel using Futures

rvest: Easily Harvest (Scrape) Web Pages

tidyverse: Easily Install and Load the 'Tidyverse'

RmecabKo: Rcpp Wrapper for Eunjeon Project

plotly: Create Interactive Web Graphics via 'plotly.js'

wordcloud2: Create Word Cloud by 'htmlwidget'

sqldf: Manipulate R Data Frames Using SQL

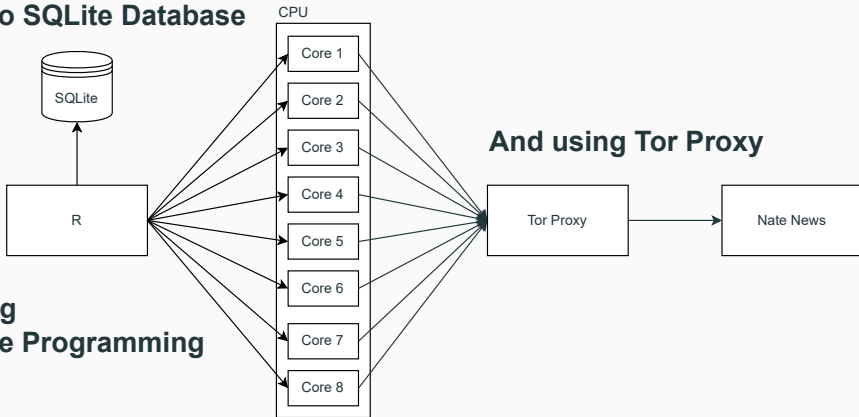
httptr: Tools for Working with URLs and HTTP

Hugging Face: Datasets, Transformers.

Pandas: Python Data Analysis Library

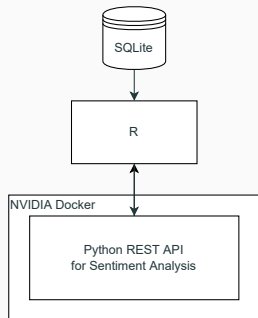
flask: A simple framework for building complex web applications.

Finally, store to SQLite Database



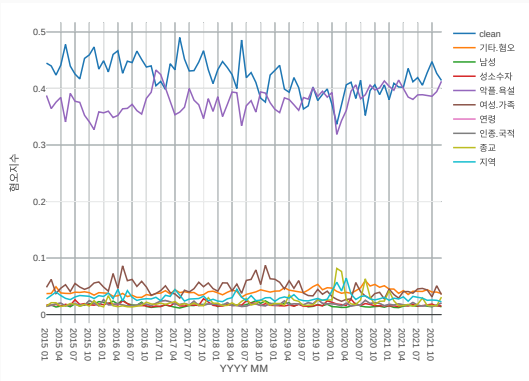
**Web Scrapping
with Multi Core Programming**

And using Tor Proxy



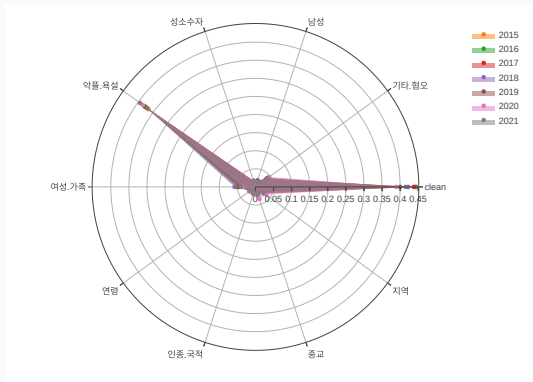
- ① SQLite로부터 수집한 댓글을 가져오기.
- ② 댓글마다 R의 `httr::POST(url)`을 이용하여 도커 내 REST API 서버에 POST 요청을 보냄.
- ③ Docker 안에서 BERT 모델을 이용해 감정 분석을 GPU로 가속하며 진행.
- ④ R에서 감정 분석 결과들을 `do.call(rbind, dfs)`를 이용해 병합.

시각화 – 꺾은선 그래프



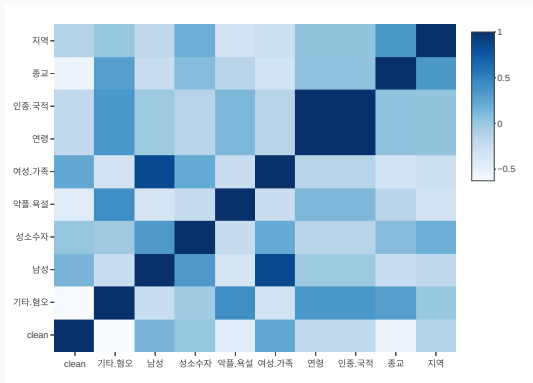
- clean과 악플의 비율은 비슷하다.
- 남성과 여성의 양상이 유사하다.
→ 상관분석을 해보자.

시각화 – 레이더 차트



- 여러 분류에 대한 분포를 쉽게 볼 수 있는 레이더 차트를 이용하여 시각화해보았다.
- clean과 악플이 다른 분류에 비해 매우 높아 보기 어려워 시연 파트에서 인터랙티브한 그래프를 보여주며 살펴보고자 한다.

시각화 – 히트 맵



- `cor(df)`를 이용하여 상관계수를 구한 결과를 히트 맵을 이용하여 시각화해보았다.
- 위에서 꺾은선 그래프를 통해 확인한 바와 같이 남성과 여성의 상관관계가 매우 높다는 사실을 알 수 있다.
- 그 외에도 지역과 종교 간의 상관성과 남성과 성소수자 간의 상관성 또한 높음을 히트 맵을 통해 확인할 수 있었다.

감정 분석

악플과 그렇지 않은 댓글의 비율은 비슷했다.

상관 분석

여성에 대한 혐오와 남성에 대한 혐오가 매우 큰 상관성을 갖는 것으로 나왔으며, 종교와 지역 간의 상관성과 남성과 성소수자 간의 상관성 또한 다른 요소에 비해 높은 것으로 나타났다.

```
1 get_news_list <- function(date, week) {
2   parameter <- sprintf("sc=&p=%s&date=%s", ifelse(week, "week", "day"), date)
3   url <- paste("https://news.nate.com/rank/cmt", parameter, sep = "?")
4   doc <- rvest::read_html(url)
5   tail <- doc > rvest::html_elements(".mduSubject > *")
6   head <- doc > rvest::html_elements(".mduSubjectList")
7   extract_tail_element <- function(e) {
8     a <- rvest::html_element(e, "a")
9     ems <- rvest::html_elements(e, "em")
10    list(
11      rank = as.integer(rvest::html_text(ems[1])),
12      comment =
13        rvest::html_text(ems[2]) >
14          stringr::str_remove_all(",") >
15            as.integer(),
16      title = rvest::html_text(a),
17      url = paste0("https:", rvest::html_attr(a, "href"))
18    )
19  }
20  extract_head_element <- function(e) {
21    ems <- rvest::html_elements(e, "em")
22    a <- rvest::html_element(e, "a")
23    list(
24      rank = as.integer(rvest::html_text(ems[1])),
25      comment =
26        rvest::html_text(ems[2]) >
27          stringr::str_remove_all(",") >
28            as.integer(),
29      title =
30        rvest::html_element(a, "strong") >
31          rvest::html_text(),
32      url = paste0("https:", rvest::html_attr(a, "href"))
33    )
34  }
35  tail <- Map(extract_tail_element, tail)
36  head <- Map(extract_head_element, head)
37  do.call(rbind.data.frame, append(head, tail))
38 }
```

```
1 get_news_comments <- function(url, page) {
2   artc_sq <- urltools::url_parse(url)$path > substring(6)
3   mid <- urltools::param_get(url, "mid")[[1]]
4   comment_url <- paste0(
5     "https://comm.news.nate.com/Comment/ArticleComment/List?",
6     "artc_sq=%s&",
7     "order=&cmtr_fl=0&prebest=0&clean_idx=&user_nm=&fold=&",
8     "mid=%s&domain=&argList=0&best=0&return_sq=&connectAuth=N&page=%d"
9   ) >
10   sprintf(artc_sq, mid, page)
11   doc <- rvest::read_html(comment_url)
12   cmt_items <- doc > rvest::html_elements(".cmt_item")
13   do.call(rbind.data.frame, Map(function(item) {
14     list(
15       name = item >
16         rvest::html_element(".nameui") >
17         rvest::html_text() >
18         stringr::str_trim(),
19       date = item >
20         rvest::html_element(".date") >
21         rvest::html_text() >
22         substring(3) >
23         paste(substring(artc_sq, 1, 4)),
24       text = item >
25         rvest::html_element(".usertxt") >
26         rvest::html_text() >
27         stringr::str_trim()
28     )
29   }, cmt_items))
30 }
```

```
1 dbGetQuery(conn, "  
2 SELECT  
3   DISTINCT (name || date || comments.text) as uniquecol,  
4   date,  
5   ct.*  
6 FROM comments  
7 INNER JOIN classified_texts ct  
8 ON comments.text = ct.text") >  
9   mutate(ym = paste(  
10     substring(date, 13, 16),  
11     substring(date, 1, 2))) >  
12   mutate(year = as.integer(substring(date, 13, 16))) >  
13   select(!uniquecol) → joined_df  
14 data.frame(word = future_Map(  
15   nouns,  
16   sqldf("  
17 SELECT text FROM joined_df  
18 WHERE ym LIKE '2015 _' )$text) >  
19   unlist() >  
20   as.vector()) → words_df  
21 sqldf("  
22 SELECT  
23   word,  
24   COUNT(1) as count  
25 FROM words_df  
26 GROUP By word  
27 HAVING length(word) > 2  
28   AND count > 50  
29 ") > wordcloud2(rotateRatio = 0)
```

SQL를 이용하여 데이터를 가공해 데이터 프레임으로 가져와 tidyverse에 있는 여러 함수를 이용해 필요에 맞게 다시 한 번 가공한 후 RmecabKo::nouns 함수를 이용하여 명사를 추출해 워드 클라우드를 그리는 코드.

프로젝트를 진행하며 작성한 모든 코드는 아래 저장소에 전부 업로드 되어 있습니다.

<https://github.com/HyunP-dev/Social-Science-Study-using-AI-Final-Project/>

– Demonstration –