

# <빅데이터 미니프로젝트 보고서>

## 유튜브 크롤링

휴먼지능정보공학과 201710799 이 지민

### 1. 서론

-다양하고, 형식을 갖추지 않은 자유로운 데이터들이 많은 유튜브를 크롤링하는 것을 주제로 정했다.  
웹상에서 raw데이터를 수집하고, 정제하고, 분석하고, 시각화하는 전반적인 과정을 진행하였다.

### 2. 프로젝트의 내용, 방법

-프로젝트 기술서

what

먹방 유튜버들의 영상정보를 웹 크롤링하여, 유튜버들이 많이 먹는 음식종류를 알아본다.  
그리고 데이터를 순위별로 알아보기 쉽게 시각화한다.

Who

신메뉴, 유행음식을 파악해 먹방 유튜버와 먹방시청자에게 마케팅을 하는 것에 도움을 준다.

Why

먹방은 신메뉴나 유행에 예민하기 때문에 사람들에게 관심있는 음식의 종류를 알 수 있다.

How

1) 50명의 먹방 유튜버의 영상정보를 selenium, bratifulSoup 모듈을 이용하여 유튜브의 영상목록 url을 이용하여 크롤링을 한다. 영상정보를 크롤링할 때, 웹상에서 스크롤이 자동으로 200번이 내려가며 모든 영상의 정보가 저장된다. 영상의 정보 중에서 영상제목, 유튜버이름, 조회수 등을 txt파일로 저장한다.

<input type="checkbox"/> youtube crawling Ae Jung.ipynb	하루 전	1.61 MB
<input type="checkbox"/> youtube crawling amatta.ipynb	하루 전	1.38 MB
<input type="checkbox"/> youtube crawling army.ipynb	하루 전	65 kB
<input type="checkbox"/> youtube crawling Boki.ipynb	하루 전	108 kB
<input type="checkbox"/> youtube crawling Bonggil.ipynb	하루 전	77.2 kB
<input type="checkbox"/> youtube crawling changbae.ipynb	하루 전	4.7 MB
<input type="checkbox"/> youtube crawling chiyeon.ipynb	하루 전	1.94 MB
<input type="checkbox"/> youtube crawling ddeonggae.ipynb	하루 전	8.41 MB
<input type="checkbox"/> youtube crawling Doram.ipynb	하루 전	3.93 MB
<input type="checkbox"/> youtube crawling dorothy.ipynb	하루 전	6.09 MB
<input type="checkbox"/> youtube crawling ealing.ipynb	하루 전	3.11 MB
<input type="checkbox"/> youtube crawling flowerpig.ipynb	하루 전	5.2 MB
<input type="checkbox"/> youtube crawling fran.ipynb	하루 전	8.80 MB
<input type="checkbox"/> youtube crawling Fume.ipynb	하루 전	1.6 MB
<input type="checkbox"/> youtube crawling G-IN.ipynb	하루 전	1.6 MB
<input type="checkbox"/> youtube crawling Gama.ipynb	하루 전	4.12 MB
<input type="checkbox"/> youtube crawling GONGSAM.ipynb	하루 전	1.54 MB
<input type="checkbox"/> youtube crawling Gyun.ipynb	하루 전	2.41 MB
<input type="checkbox"/> youtube crawling hanzy.ipynb	한 시간 전	141 kB
<input type="checkbox"/> youtube crawling hongsound.ipynb	하루 전	6 MB
<input type="checkbox"/> youtube crawling Hongyu.ipynb	하루 전	1.45 MB
<input type="checkbox"/> youtube crawling HwaRong.ipynb	하루 전	940 kB
<input type="checkbox"/> youtube crawling Hyuji.ipynb	하루 전	1.13 MB

## crawling

```
In [1]: from selenium import webdriver
from selenium.webdriver.common.keys import Keys
from bs4 import BeautifulSoup
import time

...

driver = webdriver.Chrome('C:/Users/jwlee/Code/chromedriver/chromedriver.exe')
driver.get('https://www.youtube.com/channel/UC-8sa2ivAG8q7u5SPtPQFYA/videos')

#로그인버튼
num_of_pageDowns=200
body = driver.find_element_by_xpath('//*[@html/body']

#자음으로 스크롤을 내려서 데이터가 로드 될때까지 내려감.
while num_of_pageDowns:
    body.send_keys(Keys.PAGE_DOWN)
    time.sleep(0.3)
    num_of_pageDowns -= 1
    try:
        driver.find_element_by_xpath('//*[@id="feed-main-what_to_watch"]/button").click()
    except:
        None

page = driver.page_source
page
```

```
In [4]: from bs4 import BeautifulSoup
#결과, 클래스
soup = BeautifulSoup(page, 'html')
all_title = soup.find_all('a', 'yt-simple-endpoint style=scope ytd-grid-video-renderer')
title = [soup.find_all('a', 'yt-simple-endpoint style=scope ytd-grid-video-renderer')[n].string for n in range(0, len(all_title))]
title

Out[4]: ('[리얼먹방] 굿네는 치팅행차는 치트키.....라임인정?? (ft. 허니월로파자 소주) | Chicken with Spicy powder | REAL SOUND | ASMR MUBANG |',
'[리얼먹방] 보쌈먹방!! 이 영상보고 참으면 다이어트 외장 인정 (ft. 성글 소주) | Bossam & Oysters | REAL SOUND | ASMR MUBANG | EATING SHOW |',
'[리얼먹방] 밥보디같은 꼬막!! 꼬막비빔밥 먹방 (ft. 정국장) | Spicy Blood cockle bibimbab | REAL SOUND | ASMR MUBANG | EATING SHOW |',
'[리얼먹방] 오늘은 조물조물하게 햄버거로 한끼 때울게요 (ft. 햄버거+French fries) | Hamburger+French fries | REAL SOUND | ASMR MUBANG | EATING SHOW |',
'[리얼먹방] 떡볶이 실컷해서 똥통 시켰어요 (ft. 80개 고추장똥, 라이트똥똥볶음면) | Spicy Chicken | REAL SOUND | ASMR MUBANG | EATING SHOW |',
'[리얼먹방] 햄치즈(★)편의점 아니고.....편의점(?) 먹방!! | Korea Convenience Store Food | REAL SOUND | ASMR MUBANG |',
'[리얼먹방] 대은 푸꾸미생강살 먹방!! (★)마무리는 날치알 볶음밥 | Spicy Stir-fried Octopus | REAL SOUND | ASMR MUBANG | EATING SHOW |',
'[리얼먹방] 조개탕에 혼술★ 마무리는 당면해 칼국수!! | Steamed clams | REAL SOUND | ASMR MUBANG | EATING SHOW |',
'[리얼먹방] 도가니탕(수육 날치)가 팔팔할때 국밥이 최고!!!! (ft. 작두기) | REAL SOUND | ASMR MUBANG | EATING SHOW |',
'[리얼먹방] 매콤한 팔면에 치즈토카스 !! 궁국의 조합 | jjolmyeon(spicy noodles) | REAL SOUND | ASMR MUBANG | EATING SHOW |',
'[리얼먹방] 콩나물 국밥!! 얼큰~한 할랑 먹방!! | Spicy Altang(FishmealPyeonGae) | REAL SOUND | ASMR MUBANG | EATING SHOW |',
'[리얼먹방] 피막(?) NoNo~년 피츠(★)해무로니가... (ft. 전주불떡피자, 진로이즈백) | Pizza&Soju | REAL SOUND | ASMR MUBANG | EATING SHOW |',
'[리얼먹방] 송이버섯 (★) 파장라면 ..... (ft. 군만두, 파란치) | Spicy jajiang Paman | ASMR MUBANG | EATING SHOW |',
'[리얼먹방] 대창이 이장도는 물러가줘야 대창살밥 어날?? (ft. 물맑은소스, 막창)★치즈든 쿡자 (ft. Grilled Cow intestine | REAL SOUND | ASMR MUBANG |',
'[리얼먹방] 입맛없을때 대패삼겹살(비빔면)★치즈든 RED커워 | Spring Rolls | REAL SOUND | ASMR MUBANG | EATING SHOW |',
'[리얼먹방] 연어장조림 연어달걀 만물어먹기★참고고추와 함께라면 문제없음 (ft. Marinated Salmon | REAL SOUND | ASMR MUBANG | EATING SHOW |',
```

```
youtube_video_list = []

x = 0 → #조회수인 index
y = 1 → #업로드 시점의 index

for i in range(0, len(all_title)):
    row = []
    row.append(title[i])
    # row.append(video_time[i].strip())
    row.append(channel)
    row.append(sub_num)
    row.append(view_num[x])
    x += 2 → #조회수만 append
    row.append(view_num[y])
    y += 2 → #업로드 시점만 append
    row.append(extract_date)
    youtube_video_list.append(row) → #2차원 list를 만들어줌

youtube_video_list
['구독자 104만명',
'조회수 266만회',
'3주 전',
'2019/11/17 20:22:40',
['리얼먹방'] 대창이 이장 도는 들어거먹아 대창일밥 아님?? (ft. 불닭 소스, 맥창)★다저트는 쿠키 Grilled Cow intestine | REAL SOUND | ASMR MUKBANG',
['헛지']Hanczy',
'구독자 104만명',
'조회수 129만회',
'3주 전',
'2019/11/17 20:22:40',
['리얼먹방'] 입맛없을땐 대패삼겹살해법만★다저트는 RED커워 | Spring Rolls | REAL SOUND | ASMR MUKBANG | EATING SHOW |',
['헛지']Hanczy',
'구독자 104만명',
'조회수 112만회',
'4주 전',
'2019/11/17 20:22:40',
['리얼먹방'] 연어장으로 연어일밥 만들어먹기★청양고추와 함께라면 문제없음 Marinated Salmon | REAL SOUND | ASMR MUKBANG | EATING SHOW |',
['헛지']Hanczy',
```

2 ) spark를 이용해 각 유튜버 별 영상정보를 txt로 저장하고, 그 txt를 불러와 정규표현식을 이용해 이모티콘, 특수문자를 제거한다. 그래도 제거되지 않은 문자 또는 음식명 추출에 방해되는 단어는 stopword로 제거한다. 제거한 후 rdd를 만들고 음식명을 key로 하여 key별 빈도수를 추출한다. 추출하고 난 결과는 50명 모두 하나의 txt(Frequency.txt)에 한번에 저장한다.

```
#결과
import re
a=list()

for i in hanczyList:
    text=re.sub('/', ' ', i)
    a.append(text)
    #print(text)
print(a)

#한글만 남기기
b=[]
for j in a:
    text=re.compile('[^ㄱ-ㅎ|가-힣]*')
    result = text.sub(' ', j)
    b.append(result)
print(b)
print(type(b))
# text=re.sub('[^가3131-#J316Muc00-#uad]*', ' ', j)
```

```
7 20:22:40', '', '리얼먹방'] 밥보다많은 고막!! 고막버림법 먹방 (ft. 정국영) | Spicy Blood cookle bibimbab | REAL SOUND | ASMR MUKBANG | EATING SHOW |, [헛지]Hanczy, 구독자 104만명, 조회수 111만회, 3일 전, 2019 11 17 20:22:40', '', '리얼먹방'] 오늘은 조물하게 햄버거로 한끼 때들거요 🍔 | Hamburger+French fries | REAL SOUND | ASMR MUKBANG | EATING SHOW |, [헛지]Hanczy, 구독자 104만명, 조회수 84만회, 1주 전, 2019 11 17 20:22:40', '', '리얼먹방'] 특별한 심어해서 흥분 시켰어요 🍷 (ft. 60과 고추장병, 라이트불닭볶음면) | Spicy Chicken | REAL SOUND | ASMR MUKBANG | EATING SHOW |, [헛지]Hanczy, 구독자 104만명, 조회수 156만회, 1주 전, 2019 11 17 20:22:40', '', '리얼먹방'] 헛지2시!★편의점 아니고.....편의점(?) 먹방!!! | Korea Convenience Store Food | REAL SOUND | ASMR MUKBANG |, [헛지]Hanczy, 구독자 104만명, 조회수 370만회, 1주 전, 2019 11 17 20:22:40', '', '리얼먹방'] 때온 주꾸미삼겹살 먹방!★마루리는 날치알 볶음밥 | Spicy Stir-fried Octopus | REAL SOUND | ASMR MUKBANG | EATING SHOW |, [헛지]Hanczy, 구독자 104만명, 조회수 106만회, 1주 전, 2019 11 17 20:22:40', '', '리얼먹방'] 조개탕에 흡술★마루리는 당면해 갈국수!!! | Steamed clams | REAL SOUND | ASMR MUKBANG | EATING SHOW |, [헛지]Hanczy, 구독자 104만명, 조회수 156만회, 2주 전, 2019 11 17 20:22:40', '', '리얼먹방'] 도가 나랑수육 날씨가 쌀쌀할땐 국밥이 최고!!! (ft. 라꾸기) | REAL SOUND | ASMR MUKBANG | EATING SHOW |, [헛지]Hanczy, 구독자 104만명, 조회수 131만회, 2주 전, 2019 11 17 20:22:40', '', '리얼먹방'] 때온현 꿀맛에 치즈돈까스 !! 궁국의 조합 | jishoven(spicy noodles) | REAL SOUND | ASMR MUKBANG | EATING SHOW |, [헛지]Hanczy, 구독자 104만명, 조회수 193만회, 2주 전, 2019 11 17 20:22:40', '', '리얼먹방'] 콩나물 묵묵!! 일본~한 일동 먹방!! | Spicy Aitang(FishmealRoastedSoul) | REAL SOUND | ASMR MUKBANG | EATING SHOW |, [헛지]Hanczy, 구독자 104만명, 조회수 106만회, 2주 전, 2019 11 17 20:22:40', '', '리얼먹방'] 피떡(?) Noko~난 피소 배부르니까... (ft. 전주불닭짜까, 진로이즈) | PizzadSol | REAL SOUND | ASMR MUKBANG | EATING SHOW |, [헛지]Hanczy, 구독자 104만명, 조회수 193만회, 3주 전, 2019 11 17 20:22:40', '', '리얼먹방'] 송이버섯 🍄 짜장라면... FLEX ★ ft. 군만두, 짜장면 | Spicy Jijiang Ramen | ASMR MUKBANG | EATING SHOW |, [헛지]Hanczy, 구독자 104만명, 조회수 266만회, 3주 전, 2019 11 17 20:22:40', '', '리얼먹방'] 대창이 이장 도는 들어거먹아 대창일밥 아님?? (ft. 불닭 소스, 맥창)★다저트는 쿠키 Grilled Cow intestine | REAL SOUND | ASMR MUKBANG |, [헛지]Hanczy, 구독자 104만명, 조회수 129만회, 3주 전, 2019 11 17 20:22:40', '', '리얼먹방'] 입맛없을땐 대패삼겹살해법만★다저트는 RED커워 | Spring Rolls | REAL SOUND | ASMR MUKBANG | EATING SHOW |, [헛지]Hanczy, 구독자 104만명, 조회수 112만회, 4주 전, 2019 11 17 20:22:40', '', '리얼먹방'] 연어장으로 연어일밥 만들어먹기★청양고추와 함께라면 문제없음 | Marinated Salmon | REAL SOUND | ASMR MUKBANG | EATING SHOW |, [헛지]Hanczy,
```

```
#파일저장
f=open("C:\\Users\\jile\\code\\Frequency.txt", 'a', '-1', 'utf-8')
A=len(hanzgy)
print (A)
for i in range (0,A):
    f.writelines(str(hanzgy[i]))
    f.writelines(" ")
f.close()
```

[illegible]



[illegible]

Food	Frequency
신메뉴	441
치킨	424
매운	286
라면	279
떡볶이	267
중국당면	229
짜장면	223
뽕볶음면	210
주먹밥	192
삼겹살	179
망수육	176
파김치	173
짜파게티	171
치즈볼	171
케이크	153
청양고추	152
김치	149
치즈	145
밥	123
살비김치	109

only showing top 20 rows

4) 빈도수별로 데이터결과를 시각화 하기위해 Frequency.txt에서 compile,sub정규표현식으로 음식이름data만 추출해 only\_Food.txt로 저장하였다. 또한 동의어를 처리하고 싶어 replace정규표현식을 사용했다. wordcloud모듈을 사용하여 동의어 처리된 음식이름data를 replace\_Food.txt로 저장하여, txt에서 data의 빈도수를 자동 계산해 빈도수순으로 data를 시각화하였다.(Konply, pygame등의 모듈을 이용)

[illegible][illegible]

1320

7

1320

사진 - wordcloud\_food.png



### 3. 프로젝트 결과

```
tDf=resultDf.withColumn("Food",resultDf['_2'].cast("string")).drop('_2')
ResultDf=tDf.withColumn("Frequency",tDf['_1'].cast("integer")).drop('_1')
```

```
ResultDf.show()
```

Food	Frequency
신메뉴	441
치킨	424
매운	285
라면	279
떡볶이	257
중국당면	229
짜장면	223
불닭볶음면	210
주먹밥	192
삼겹살	179
탕수육	176
파김치	173
짜파게티	171
치즈볼	171
케이크	153
청양고추	152
김치	149
치즈	145
밥	123
살비김치	108

only showing top 20 rows



왼쪽은 raw데이터로 구한 음식이름별 빈도수이고, 오른쪽은 Frequency.txt의 data에서 정규표현식을 거쳐 음식이름별 빈도수를 토대로 데이터 시각화한 모습이다. 치킨, 매운, 라면, 떡볶이 등 순으로 음식빈도가 많이 나왔다. 또한 프로젝트를 하면서 느낀점은 웹상에 데이터는 정말 깨끗하지 않은 형태이고, 다양한 방법으로 여러번 정제해야 가치있는 데이터를 얻을 수 있겠다고 생각했다. 내가 생각한 방법으로 정제하고 분석하려면 data형태가 맞지 않거나, 원하지 않는 data까지 섞여있는 등 다양한 문제가 발생했다. 이런 문제를 해결하고 고민하는 과정이 유익했던 것 같다. 또한 수집과 정제와 분석을 잘한다면 무한한 가치가 있고, 다양한 의미가 숨겨져 있을 것이라는 기대를 하게 되었다.