

# Big Data Programming

AI 융합 학부  
20170403 최혜원

# INDEX

1. 주제

2. 알고리즘

3. 진행 상태

4. 개발 계획

5. Q & A

## 1. 주제

이 주제에 대한 설명

# “정치 뉴스와 연예 뉴스의 상관관계 분석”

가설 : 어떤 연예계 사건이 터졌을 때, 정치적, 사회적 문제를  
덮고자 하는 것일지도 모른다.

## 2. 알고리즘

구현 flow

If 특정 정치 키워드가 이슈가 된 시기를 파악

If 그 시기 전후로 어떤 연예 관련 keyword 그래프가 급격히 상승

Then 두 특정 키워드의 추이를 분석, 결과 도출

## 2. 알고리즘

구현 flow

If 특정 정치 키워드가 이슈가 된 시기를 파악

- 1) 키워드를 정한다. (ex. 국정원 해킹, 교과서 국정화 추진 논란, LH) (중간)
- 2) 정한 키워드를 리스트로 정리한다. ( 날짜데이터 등) (기말)

## 2. 알고리즘

구현 flow

If 어떤 연예 관련 keyword 그래프가 급격히 상승

- 1) 키워드 발생 전후기간에 연예 관련 데이터 추이를 파악한다.(중간)
- 2) 의미 있는 연예 관련 데이터를 추출한다.(기말)

## 2. 알고리즘

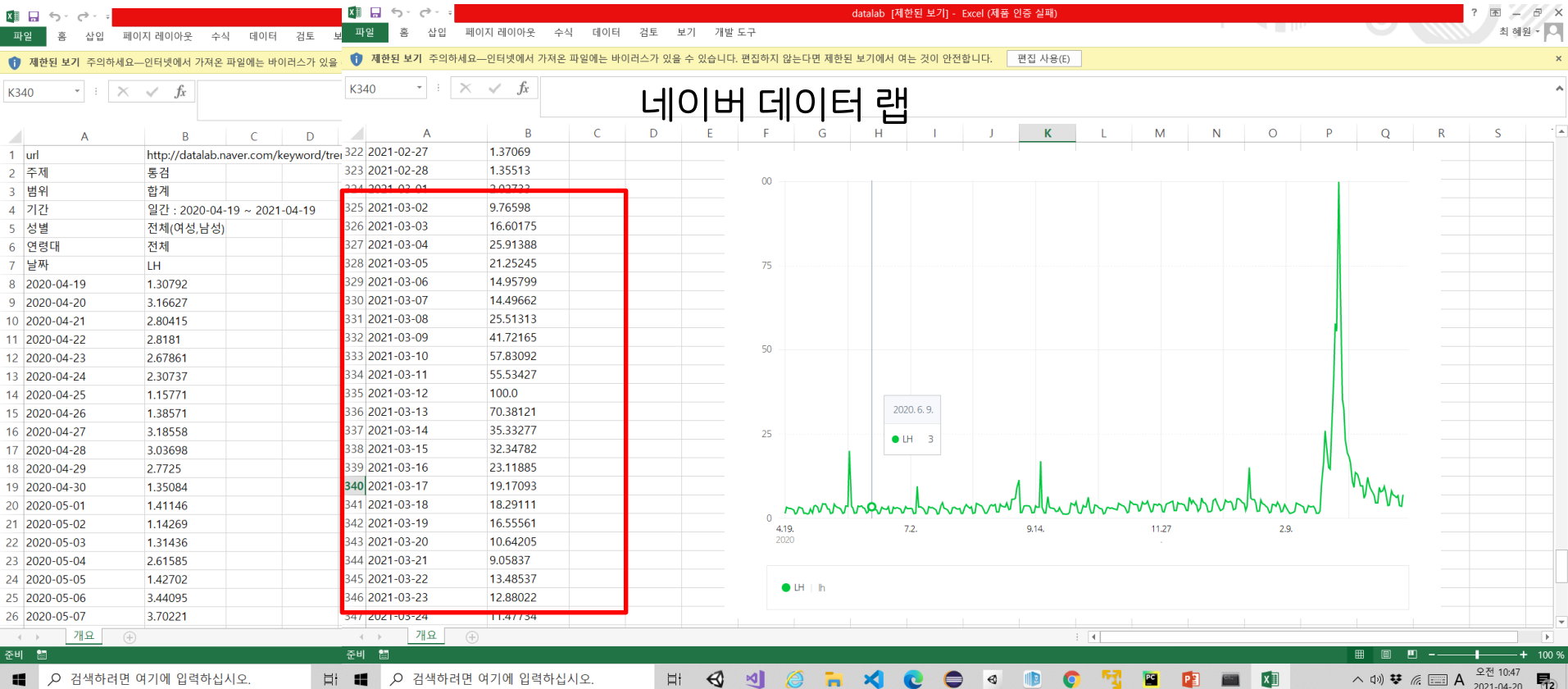
구현 flow

두 특정 키워드의 추이를 분석, 결과 도출

- 1) 정치 키워드와 연예 키워드 추이를 시각화한다.
- 2) 의미를 도출한다.

### 3.진행상태

#### 현재 프로젝트 진행 상태 - 데이터 수집





### 3. 진행 상태

#### 현재 프로젝트 진행 상태 - 데이터 수집

- 트위터의 파이썬 라이브러리인 **tweepy**를 활용하여 과거 및 현재 인기 트위터 데이터를 가져와서 가공하였다.
- Tweet volume 이 일정 수 이상 큰 경우의 단어를 모으기 위해 volume 값들을 출력하였다.
- 날짜는 cursor 메서드의 since값을 지정하여 원하는 날짜의 data를 가져올 수 있다.

파일(F) 편집(E) 선택 영역(S) 이동(G) 실행(R) 터미널(T) 도움말(H)

탐색기

▼ 열려 있는 편집기

natePann.py

natePann.json M

titleNaver.py U

twitter\_Seoul\_trend.json

available\_locs\_for\_trend.json

Ndata.json U

twitter\_Korea\_trend.json M X

▼ PYTHON

\_\_pycache\_\_

config.cpython-38.pyc M

data

locs

available\_locs\_for\_trend.js...

melon.json

natePann.json M

Ndata.json U

twitter\_Korea\_trend.js... M

twitter\_Seoul\_trend.json

twitter\_Worldwide\_trend.json

config.py

melon.py

natePann.py

titleNaver.py U

twitter.py M

data > {} twitter\_Korea\_trend.json > {} 0 > [ ] trends

318 },

319 {

320 "name": "그럴만두",

321 "url": "http://twitter.com/search?q=%EA%B7%B8%EB%9F%B4%EB%A7%8C%EB%91%90",

322 "promoted\_content": null,

323 "query": "%EA%B7%B8%EB%9F%B4%EB%A7%8C%EB%91%90",

324 "tweet\_volume": null

325 },

326 {

327 "name": "에너지 낭비",

328 "url": "http://twitter.com/search?q=%22%EC%97%90%EB%84%88%EC%A7%80+%EB%82%AD%EB%B9%84%22",

329 "promoted\_content": null,

330 "query": "%22%EC%97%90%EB%84%88%EC%A7%80+%EB%82%AD%EB%B9%84%22",

331 "tweet\_volume": null

332 },

333 {

334 "name": "taemin",

335 "url": "http://twitter.com/search?q=taemin",

336 "promoted\_content": null,

337 "query": "taemin",

338 "tweet\_volume": 183058

339 },

340 {

341 "name": "플라잉키스",

342 "url": "http://twitter.com/search?q=%ED%94%8C%EB%9D%BC%EC%9E%89%ED%82%A4%EC%8A%A4",

343 "promoted\_content": null,

344 "query": "%ED%94%8C%EB%9D%BC%EC%9E%89%ED%82%A4%EC%8A%A4",

345 "tweet\_volume": null

346 },

347 {

348 "name": "원더케이",

349 "url": "http://twitter.com/search?q=%EC%9B%90%EB%8D%94%EC%BC%80%EC%9D%B4",

350 "promoted\_content": null,

351 "query": "%EC%9B%90%EB%8D%94%EC%BC%80%EC%9D%B4",

터미널

디버그 콘솔

문제

출력

1: powershell

PS C:\Users\soohi\Desktop\bigdata\real\_time\_trend\_caht\_kakao\python> [ ]

main\* Python 3.9.0 64-bit 0 0 0

줄 304, 열 6 공백: 4 UTF-8 with BOM CRLF JSON with Comments

검색하러 면 여기에 입력하십시오.

오후 12:00 2021-04-20

### 3. 진행 상태

현재 프로젝트 진행 상태 - 데이터 수집

구글 트렌드 api(google-trends-api)를 통해서 원하는 데이터를 가져온다.  
(날짜, 키워드 등)

1위 :: 슈퍼리그  
2위 :: 비트코인  
3위 :: 이현배  
4위 :: 버팀목자금플러스  
5위 :: 무리뉴 경질  
6위 :: 기모란  
7위 :: 바르셀로나  
8위 :: 아스날  
9위 :: F1

```
data > natePann.json > ...
1  [
2  [
3  "1",
4  "이하늘 왜 화난지 물어놓음"
5  ],
6  [
7  "2",
8  "온 몸이 기름진 인간이 꼭 하는 루틴"
9  ],
10 [
11 "3",
12 "킹덤 또 논란 터졌다 .."
13 ],
14 [
15 "4",
16 "이하늘 김창열 어떻게 생각해요???"
17 ],
18 [
19 "5",
20 "우리나라 남들 억대 비주얼 댄들"
21 ],
22 [
23 "6",
24 "서울살아도 잘 안타보는 지하철 노선 갑"
25 ],
26 [
27 "7",
28 "역대 걸그룹 비주얼"
29 ],
30 [
31 "8",
32 "너네 예리 인스타 봤냐???"
33 ],
34 [
35 "9",
36 "아이유로 살기 vs 그냥 살기"
37 ],
38 ]
```

간내평균

43 %

.72 %

86 %

96 %

## 4. 개발 계획

앞으로의 개발 계획 - 데이터 가공

1. 정치적으로 굵직한 사건들을 시기별로 정리할 것.  
(선별 방법 : 검색량의 **볼륨 변화량**이 급증했을 경우 정치적 굵직한 사건이라고 판단)
2. 그 시기 전후로 연예관련 데이터를 모을 것.
3. 뉴스, 커뮤니티, 트위터 이외의 연예계 관련 데이터 소스(뉴스 기자, 뉴스 사이트: 피드백 반영)를 찾을 것.
4. 최대한 많은 데이터를 수집하여 의미 있는 데이터만 전처리하여 시각화할 것.



**QnA**

**Thank you**