



개인 포트폴리오_김민준

빅데이터를 활용한 **SNS 인기 여행지**

대한민국 구석구석 SNS 인기 여행지 데이터 수집, 저장 및 RPA 활용

대한민국 구석구석 + 인스타그램 + RPA

1. 추진 목표

SNS 인기 여행지 데이터 수집, 저장 및 RPA 활용

- 대한민국 구석구석의 인스타그램 인기 여행지 데이터 수집 및 저장
- 수집 데이터를 저장 및 RPA 활용(이메일 전송, 카카오톡 메시지 전달)



2. 과제 수행 범위

대한민국  구석구석

인스타그램 인기 여행지
TOP 11 수집

여행지명, 지역, 상세정보,
해시태그, 여행지별 이미지 3장

Instagram

여행지명에 따른 게시물
개수 수집

인기 여행지일수록 게시물
개수 증가

NAVER

네이버 날씨 크롤링

오늘,내일,모레

3. 요구사항 정의

인스타그램
게시물 수집

여행지 명이 불명확할 경우
공백처리

CSV, JSON
파일 저장

업무용 혹은 웹 데이터 활용

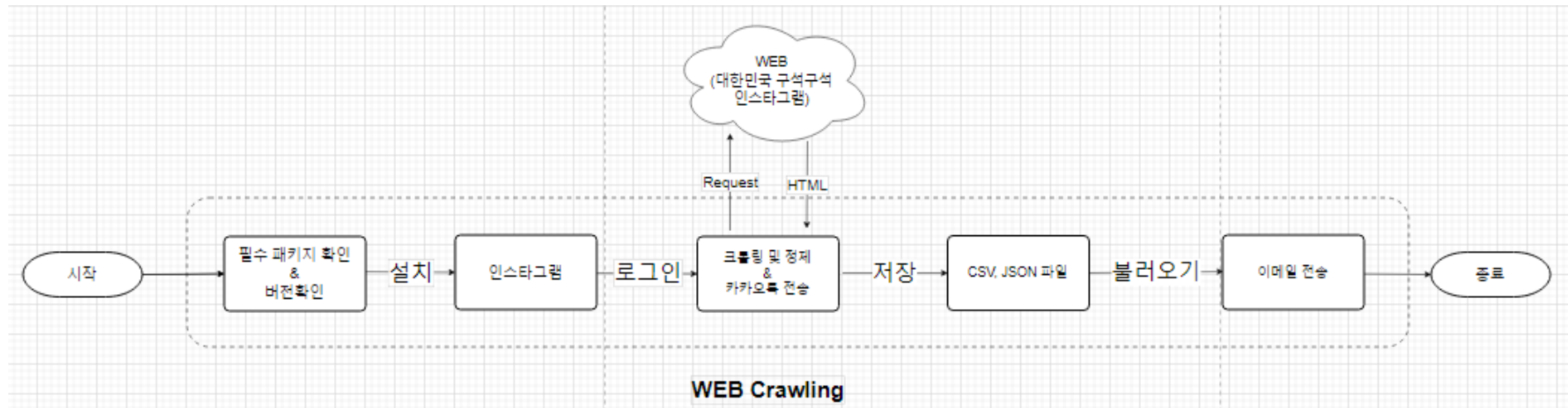
파일 이메일 전송

이메일 자동화

카카오톡 메시지
전송

카카오톡 자동화

4. 프로세스 설계



1. 필수 패키지 확인 & 버전확인

데이터 수집에 필요한 패키지 설치 및 최신버전 확인

2. 인스타그램

데이터 수집 전 인스타그램 로그인

3. 크롤링 및 정제 & 카카오톡 전송

대한민국 구석구석 수집 및 여행지 명에 맞는
인스타그램 수집, 수집 내용 카카오톡 피드
형식으로 전송

4. CSV, JSON 파일 저장

데이터 수집내용 csv 파일과 json 파일 저장

5. 이메일 전송

csv파일과 json 파일을 이메일 전송

5. 구현

5-1) 필수 패키지 확인 및 버전 확인

```
1 print('='*120)
2 print('
3 필요 모듈 설치 여부를 확인합니다.\n
4 체크 모듈 : requests, pandas,selenium, bs4
5 ')
6 print('='*120)
7
8 import sys # 파이썬 인터프리터의 변수와 함수를 직접 제어하기 위해 호출
9 import subprocess # 모듈들 확인 및 설치를 위한 셸 명령 실행 모듈 호출
10
11 try :
12     import requests # 웹페이지 url 처치를 위한 requests 모듈 호출
13     import pandas # 데이터프레임 저장을 위한 pandas 모듈 호출
14     import selenium # 웹페이지 제어를 위한 selenium 모듈 호출
15     import bs4 # html파싱을 위한 BeautifulSoup 모듈 호출
16
17
18 except :
19     # pip 모듈 업그레이드
20     subprocess.check_call([sys.executable, "-m", "pip", "install", "--upgrade", "pip"])
21     # requests 모듈 최신버전 설치
22     subprocess.check_call([sys.executable, "-m", "pip", "install", "--upgrade", "requests"])
23     # pandas 모듈 최신버전 설치
24     subprocess.check_call([sys.executable, "-m", "pip", "install", "--upgrade", "pandas"])
25     # selenium 모듈 최신버전 설치
26     subprocess.check_call([sys.executable, "-m", "pip", "install", "--upgrade", "selenium"])
27     # BeautifulSoup 모듈 최신버전 설치
28     subprocess.check_call([sys.executable, "-m", "pip", "install", "--upgrade", "bs4"])
29     print('='*120)
30     print('필요 모듈 설치가 완료되었습니다.')
31     print('='*120)
32
33 else:
34     print('='*120)
35     print('필요 모듈 설치가 완료되어 버전을 확인합니다.')
36     print('='*120)
37     # BeautifulSoup 모듈 최신버전 확인
38     print('requests module version check :'+ bs4.__version__)
39     # pandas 모듈 최신버전 확인
40     print('pandas module version check :'+ pandas.__version__)
41     # selenium 모듈 최신버전 확인
42     print('selenium module version check :'+ selenium.__version__)
43     # requests 모듈 최신버전 확인
44     print('requests module version check :'+ requests.__version__)
45     print('='*120)
```

=====

필요 모듈 설치 여부를 확인합니다.

체크 모듈 : requests, pandas,selenium, bs4

=====

=====

필요 모듈 설치가 완료되어 버전을 확인합니다.

=====

requests module version check :4.10.0
pandas module version check :1.3.4
selenium module version check :4.1.0
requests module version check :2.26.0

=====

5. 구현

5-2) 인스타그램

```
1 '''
2 함수명 : instagram_login
3 기능 : 인스타그램 로그인(id,pw) 입력
4 매개변수 : user_id , user_pw
5 - user_id : (str) 인스타그램 아이디
6 - user_pw : (str) 인스타그램 비밀번호
7 리턴값 : 없음
8 '''
9
10 # 인스타그램 로그인 함수 정의
11 def instagram_login(user_id, user_pw) :
12     # 아이디 입력
13     id_input = driver.find_element_by_xpath('//*[@id="loginForm"]/div/div[1]/div/label/input') # 인스타그램 로그인 입력창
14     id_input.send_keys(user_id) # 인스타그램 아이디 입력
15     time.sleep(1) # 아이디 입력 후 1초 대기
16
17     # 비밀번호 입력
18     pw_input = driver.find_element_by_xpath('//*[@id="loginForm"]/div/div[2]/div/label/input') # 인스타그램 비밀번호 입력창
19     pw_input.send_keys(user_pw) # 인스타그램 비밀번호 입력
20     time.sleep(1) # 비밀번호 입력 후 1초 대기
21     pw_input.send_keys(keys.ENTER) # 비밀번호 입력창에 엔터 입력
22     time.sleep(120)
23
24
25 '''
26 함수명 : wake_cookie
27 기능 : 인스타그램 로그인 성공 쿠키를 받음
28 매개변수 : driver , user_id
29 - driver : 웹 드라이버 객체 호출
30 - user_id : (str) 인스타그램 아이디
31 리턴값 : 없음
32 '''
33
34 def wake_cookie(driver, user_id) : # 쿠키받기기 함수정의
35     with open(f'./cookies/instagram_{user_id}_cookies.pkl', 'wb') as f : # cookies폴더에 인스타그램 아이디와 동일한 pickle(쿠키) 파일 만들기
36         pickle.dump(driver.get_cookies(),f) # 웹드라이버로 쿠키를 가져온다.
37
38
39 '''
40 함수명 : load_cookie
41 기능 : 인스타그램 로그인 성공 쿠키를 불러오기
42 매개변수 : driver , user_id
43 - driver : 웹 드라이버 객체 호출
44 - user_id : (str) 인스타그램 아이디
45 리턴값 : 없음
46 '''
47
48 def load_cookie(driver, user_id) : # 쿠키 불러오기 함수정의
49     with open(f'./cookies/instagram_{user_id}_cookies.pkl', 'rb') as f : # cookies폴더에 인스타그램 아이디와 동일한 pickle(쿠키) 파일 불러오기
50         cookies = pickle.load(f) # cookies 변수에 쿠키값들을 불러오고, 저장
51
52         for cookie in cookies : # cookies 만큼 cookie를 반복
53             driver.add_cookie(cookie) # 웹드라이버에 cookie 정보를 추가
54
55
56
57 # 쿠키받기기 로그인
58 if os.path.isfile(f'./cookies/instagram_{user_id}_cookies.pkl') : # 경로에 쿠키파일이 있다면
59     print('='*120)
60     print('쿠키가 존재합니다. 인스타그램 로그인에 성공하였습니다.')
61     print('='*120)
62     load_cookie(driver, user_id) # 쿠키 불러오기 함수 호출
63     driver.refresh() # 웹드라이버 새로고침
64
65
66 else : # 경로에 쿠키파일이 없다면
67     print('='*120)
68     print('쿠키가 존재하지 않습니다. 인스타그램 로그인을 시도합니다.')
69     print('='*120)
70     instagram_login(user_id, user_pw) # 인스타그램 로그인 함수 호출
71     wake_cookie(driver, user_id) # 쿠키 받기기 함수 호출
```

=====

쿠키가 존재합니다. 인스타그램 로그인에 성공하였습니다.

=====

5. 구현

5-3) 크롤링 및 정제 & 카카오톡 전송

```
1 '''
2 함수명 : get_korea
3 기능 : 대한민국 구석구석 SNS 인기여행지 TOP11 수집
4 매개변수 : 없음
5 리턴값 : 없음
6
7 '''
8
9 def get_korea() : # 대한민국 구석구석 top 11 수집 함수 정의
10     #리턴 값 고민하기
11     for x in range(1,7) : # 1부터 6까지 반복
12         driver.find_element_by_xpath(f"/html/body/div[3]/div[3]/div[1]/div[1]/ul[1]/li[{x}]").click() # 대한민국 구석구석 인스타그램 추천 여행지 1번째
13         time.sleep(3)
14         detail_page() # 상세페이지 함수 호출
15         time.sleep(3) # 페이지 로딩을 위한 3초 대기
16
17
18
19 # 넘어 갈 때 자바 스크립트 불러 여러 발생
20 for y in range(1,6) : # 1부터 5까지 반복
21     driver.find_element_by_xpath(f"/html/body/div[3]/div[3]/div[1]/div[1]/ul[2]/li[{y}]").click() # 대한민국 구석구석 인스타그램 추천 여행지 2번째
22     time.sleep(3) # 페이지 로딩을 위한 3초 대기
23     detail_page() # 상세페이지 함수 호출
24
25
26
27 save_data.csv_save() # save_data 모듈의 csv 저장 함수호출
28 time.sleep(2) # 로딩을 위한 2초 부여
29 save_data.json_save() # save_data 모듈의 json 저장 함수호출
30 time.sleep(2) # 로딩을 위한 2초 부여
31 send_mail.email(full_dir) # send_mail 모듈의 email 함수 호출(csv,json파일을 이메일로 보내준다.)
32
33 driver.close() # 창닫기
```

=====

설악산국립공원 사진 저장을 완료하였습니다.

=====

여행지 : 설악산국립공원 / 강원 속초시

검색태그 : 설악산국립공원 / 인스타그램 게시물 수 : 44,097 건

상세정보 : 설악산은 강원도 속초시를 포함해 고성과 인제, 양양군에 걸쳐 있는 유명한 자연 명소이다. 1970년 한려해상국립공원에 이어 다섯 번째로 국립공원 명칭을 받았으며, 1982년에는 국내 최초로 유네스코 생물권보전지역에 지정됐다. 또한 국제적으로 우수하게 관리되고 있는 보호지역 명단인 세계자연보전연맹(IUCN) 녹색목록에도 등재되어 있다. 설악산국립공원은 총면적이 약 398km²에 이를 만큼 광대한 규모를 자랑한다. 주봉인 대청봉을 기준 삼아 속초시에 접한 동쪽은 외설악, 서쪽은 내설악, 한계령과 오석 방면은 남설악으로 구분한다. 한라산과 지리산 다음으로 높은 대청봉은 1년 중 5~6개월은 눈에 덮여 있는데 ‘설악’이란 이름은 여기서 유래되었다. 대청봉에 오르면 수려한 자연 경관을 품은 설악산 전경과 동해 바다를 조망할 수 있으며 특히 일출과 낙조가 일품으로 꼽힌다. 대청봉 외에도 소청봉, 화채봉, 중청봉 등 30여개의 산봉우리와 웅장한 장관을 이룬다. 이외에 비룡폭포, 울산바위, 흔들바위, 금강굴 등 숨은 비경들이 가득하다. 설악산 소공원 인근에 운영 중인 케이블카를 이용하면 노약자나 장애인도 어렵지 않게 설악산 풍경을 감상할 수 있다.[체험프로그램]등산, 트레킹

URL : https://korean.visitkorea.or.kr/detail/ms_detail.do?cotid=4db10875-c210-476a-9873-dc307faccf47&big_category=A01&mid_category=A0101&big_area=32

해시태그 : #21-22한국관광100선#21_22한국관광100선#강원권#강원도_여행지_추천#관광지#국립공원#다양한식물#비룡폭포#설악종아영장#설악산국립공원#식물자원외_보고

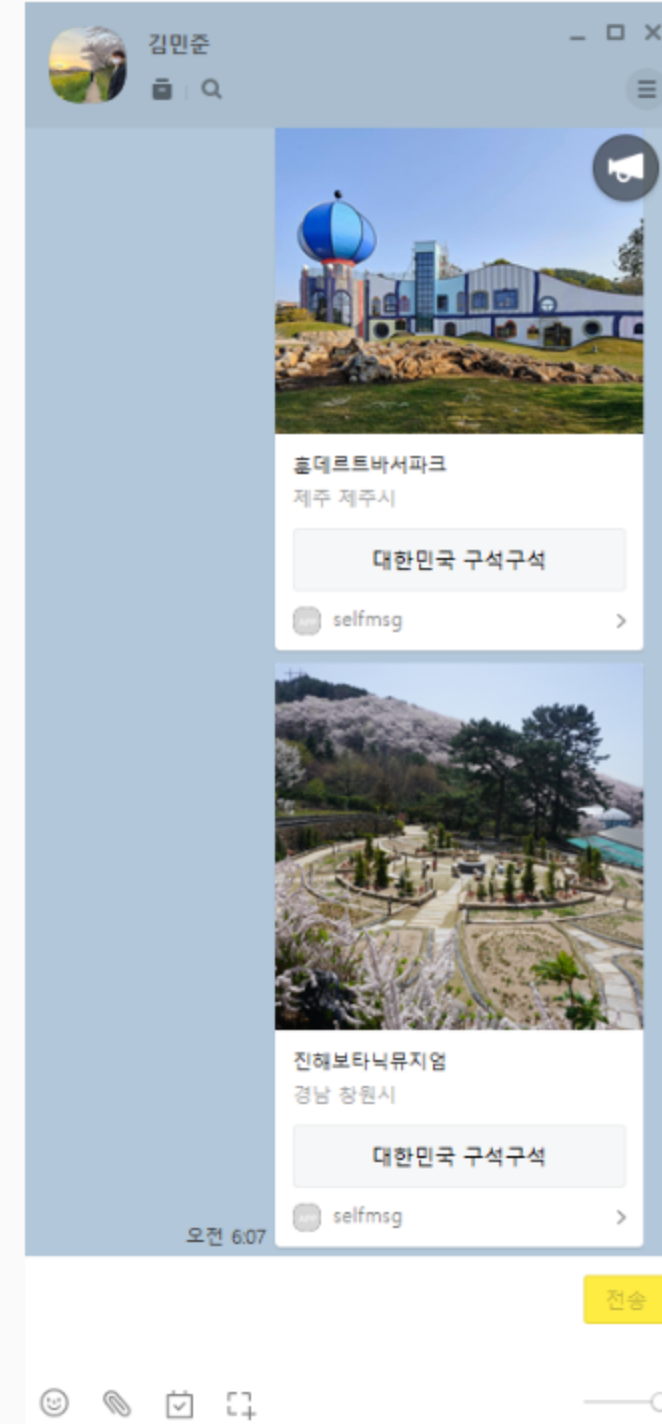
=====

카카오톡 메시지를 보냅니다.

=====

카카오톡 메시지 전송 성공하였습니다.

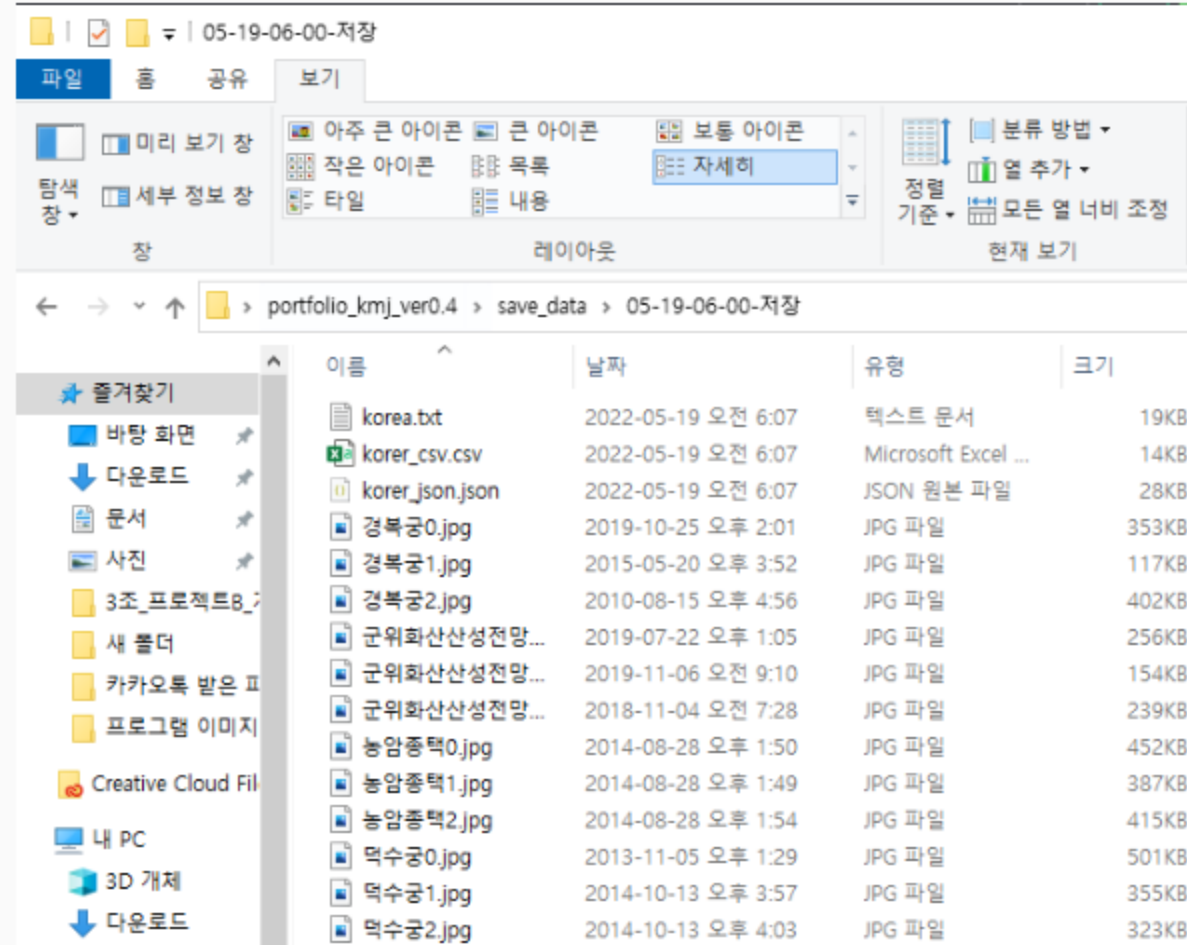
=====



5. 구현

5-4) CSV, JSON 저장

```
1 '''
2 함수명 : csv_save
3 기능 : txt 파일을 csv 파일로 저장하는 함수
4 매개변수 : 없음
5 리턴값 : 없음
6 '''
7 def csv_save() : # txt 파일 -> csv 파일로 저장 함수 정의
8
9     print('='*120)
10    print('CSV 파일로 저장을 시작합니다') # csv 저장 시작 출력
11    print('='*120)
12    try : # 해당 경로에 txt파일이 있다면 파일 불러오기
13        csv_data = pd.read_csv('./korea.txt', sep='\t', encoding='utf-8', lineterminator='\n', error_bad_lines=False)
14    except : # 해당 경로에 txt파일이 없을 경우
15        print('='*120)
16        print('CSV 파일로 저장을 실패했습니다.(파일유무를 확인해주세요.)')
17        print('='*120)
18
19    else :
20        # csv 실행
21        csv_data.columns = ['여행지명', '위치', '인스타그램 게시물 수', '상세정보', '해시태그']
22        # csv 파일로 저장
23        csv_data.to_csv('./korer_csv.csv', encoding='utf-8-sig', index=False)
24        # 저장 완료 출력
25        print('='*120)
26        print('CSV 파일로 저장을 완료했습니다.')
27        print('='*120)
28
29 #csv_save(full_dir) # 테스트 출력
30
31 '''
32
33 함수명 : json_save
34 기능 : csv 파일을 json 파일로 저장하는 함수
35 매개변수 : 없음
36 리턴값 : 없음
37
38 '''
39 def json_save() : # csv파일 -> json 파일로 저장 함수 정의
40
41     print('='*120)
42     print('JSON 파일로 저장을 시작합니다') # json 저장 시작 출력
43     print('='*120)
44     try : # 해당 경로에 csv파일이 있다면 파일 불러오기
45         json_data = pd.read_csv('./korer_csv.csv', sep=',', encoding='utf-8')
46     except : # 해당 경로에 csv파일이 없을 경우
47         print('='*120)
48         print('JSON 파일로 저장을 실패했습니다.(csv 파일유무를 확인해주세요.)')
49         print('='*120)
50
51     else :
52         # json 파일로 저장
53         json_data.to_json('./korer_json.json', orient='records')
54         # 저장완료 출력
55         print('='*120)
56         print('JSON 파일로 저장을 완료했습니다.')
57         print('='*120)
```



```
=====
CSV 파일로 저장을 시작합니다
=====
c:\Users\jmk\Desktop\portfolio_kmj_ver0.4_20220520\portfolio\complete_code\korea_crawling.py:180: FutureWarning: The error
and will be removed in a future version.

save_data.csv_save() # save_data 모듈의 csv 저장 함수호출
=====
CSV 파일로 저장을 완료했습니다.
=====
JSON 파일로 저장을 시작합니다
=====
JSON 파일로 저장을 완료했습니다.
=====
```

5. 구현

5-5) 이메일 전송

```
1  ....
2  함수명 : email
3  매개변수 : full_dir
4  - full_dir : 파일 저장된 경로가 저장된 변수
5  리턴값 : 없음
6
7  ...
8  def email(full_dir): # 이메일 보내는 함수 정의
9      print('='*120)
10     print('메일을 보내는 중입니다.') # 이메일 시작 메시지 출력
11     print('='*120)
12
13     # 전자메일 메시지 작성을 위한 모듈 호출
14     msg = EmailMessage()
15     msg["Subject"] = "대한민국 구석구석 크롤링 파일입니다." # 이메일 제목
16     msg["From"] = email_address # 보내는 사람
17     msg["To"] = "tkffwnj2002@naver.com" # 받는 사람
18     msg.set_content('대한민국 구석구석 크롤링한 파일데이터입니다.') # 이메일 본문 텍스트
19
20     # 구글 -> mimetype을 검색하면 검색하는 포맷이 나온다.
21     # 보내려는 파일 타입을 찾은 후 메인,서브 속성에 넣으면 된다.
22     # json => minatype='application' , sub = json
23     # csv => text , csv
24     # 파일이름에 문자열을 입력하면 입력값으로 파일이 들어간다.
25
26     # csv파일 불러오기
27     with open('./korer_csv.csv', 'rb') as f_csv :
28         # msg 객체에 파일 첨부할 객체화(open한 파일을 읽고, 메인타입=, 서브타입=, 파일이름=)
29         msg.add_attachment(f_csv.read(), maintype='text', subtype='csv', filename = 'korer_csv')
30
31     # json파일 불러오기
32     with open('./korer_json.json', 'rb') as f_json :
33         # msg 객체에 파일 첨부할 객체화(open한 파일을 읽고, 메인타입=, 서브타입=, 파일이름=)
34         msg.add_attachment(f_json.read(), maintype='application', subtype='json', filename = 'korer_json')
35
36     # 이메일 보내기
37     with smtplib.SMTP("smtp.gmail.com",587) as smtp :
38         smtp.ehlo() # 연결이 잘 되는지 확인
39         smtp.starttls() # 메일 내용이 암호화되어 전송
40         smtp.login(email_address,email_pw) # 로그인에 필요한 이메일 주소 및 비밀번호 함수 호출
41         smtp.send_message(msg) # 메시지 보내기
42
43
44     print('='*120)
45     print('메일을 보내기를 완료했습니다.') # 이메일 완료 메시지 출력
46     print('='*120)
```

☆ 대한민국 구석구석 크롤링 파일입니다. [🔗](#)

2022-05-21 (토) 01:52

보낸사람 VIP <tkffwnj96@gmail.com>

받는사람 <tkffwnj2002@naver.com>

📎 일반 첨부파일 2개 (41KB) 모두 저장

🛡️ 파일 저장 시 바이러스 검사 자동 수행

📎 korer_csv 14KB

✖

📎 korer_json 27KB

✖

대한민국 구석구석 크롤링한 파일데이터입니다.

=====

메일을 보내는 중입니다.

=====

메일을 보내기를 완료했습니다.

=====

THANK YOU



이름 / 연락처

김민준 / 010-9608-2015



이메일

tkffuwnj2002@naver.com



GitHub

<https://github.com/jmk2015>
git/portfolio_koreatrip.git