

[Python] BeautifulSoup을 통한 이미지 스크래핑 하기 — 나무늘보의 개발 블로그

노트북: 첫 번째 노트북

만든 날짜: 2020-10-25 오후 5:14

URL: <https://continuous-development.tistory.com/108?category=736681>

Python

[Python] BeautifulSoup을 통한 이미지 스크래핑 하기

2020. 10. 7. 17:38 수정 삭제 공개

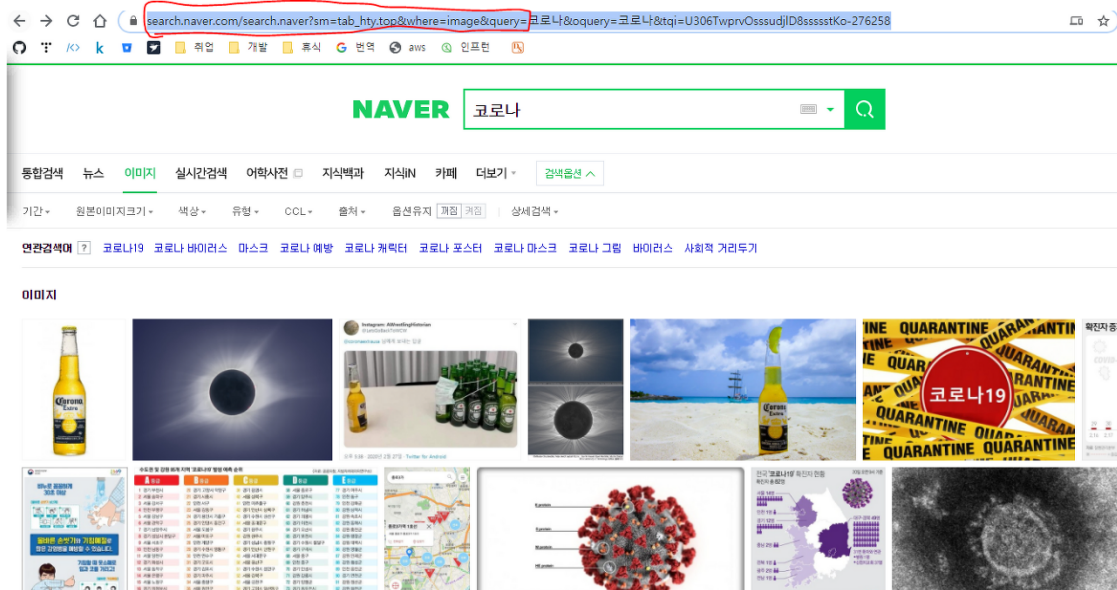
※환경은 anacond + jupyter에서 하였습니다.

```
In [3]: from urllib.request import urlopen
        from bs4 import BeautifulSoup
        from urllib.error import HTTPError
        from urllib.error import URLError
        import pandas as pd
        from urllib.parse import quote_plus
```

크롤링을 하는 데 있어서 기본적으로 필요한 함수들이다. 없으면 anaconda prompt에서 pip install을 통해 받도록 하자

내가 크롤링하려는 사이트는 네이버다. 네이버에서 검색어를 입력하고 그 이미지를 가져오기로 하자.

첫 번째로 url을 가져오자.



저렇게 검색어를 기준으로 앞에 url을 가져온다. 검색어를 입력하는 것에 따라 이미지를 가져오기 위해서이다.

base_url에 검색어 url을 넣어주고

keyword에는 내가 넣을 검색어를 넣어주기 위해 input을 통해 검색어를 입력한다.

ImgCnt는 스크랩할 이미지 개수를 정해준다.

```
In [5]: base_url = 'https://search.naver.com/search.naver?sm=t&where=image&query=' # base url 이다. 여기서 뒤에 어떤 검색어를 넣냐에 따라 달라진다.
keyword = input('검색어 입력 : ') # 여기서 네이버에서 이미지 검색하던것 처럼 명령어를 입력해준다.
imgCnt = int(input('스크랩 할 이미지 개수 : ')) # 스크랩 할 이미지 개수를 지정해준다.
try:
    html = urlopen(base_url+quote_plus(keyword)) # quote_plus 를 통해 URL Encoding을 한다.
except HTTPError as he:
    print('http error')
except URLError as us:
    print('url error')
else:
    soup = BeautifulSoup(html.read(), 'html.parser')

검색어 입력 : 코로나
스크랩 할 이미지 개수 : 50
```

위와 같은 식으로 로직을 만들어준다. 내가 검색어를 입력하는 거에 따라 결과가 달라지는 크롤링 코드를 만들어보자.

사이트를 스크래핑하고 내가 원하는 부분이 어디 있는지 확인한다.

https://search.static.net/common/?src=http%3A%2F%2Fpost.phinf.naver.net%2Fimage%2FMDA0TNTMTgZGFMAxQzN2NmNmM0ODM4.petATgz0Xbomp4_33Sj0ze2VvL26Y8hc959QxUPBq.f83o-XtUuyCld7wGsnCuoh7rURFCB55wBUl3m1.Qq.JPEGiZFJfhtAsVkvPsvUDeZOsnk3Wc3IqJ.pjg&type=b400

https://search.static.net/common/?src=http%3A%2F%2Fcafefiles.naver.net%2Fimage%2FYdAYMdyIMyJfOTQZGFMDAxtNGVwZjU0NTg5M2NmNmM0UHUY9Uw4qeGc-tnsCSHw.xfc9GnkMsog7oa3fg.jzVlcsmHYVMPC_NrnNMVE2n_2IdAnQYoIKto8b01Oq.FNgiK2F2Fts50dX25B7Z25CE5E3B32gmA.X2FB8BX25E7425C1X25F8.ngnqttype=b400

https://search.static.net/common/?src=http%3A%2F%2Fimages.naver.net%2Fimage%2F001X2F2020X2F02X2F292F2FAK2R202022901000087_01_i_20200229051206365_pjg&type=b400

https://search.static.net/common/?src=http%3A%2F%2Fblogfiles.naver.net%2F0100126_285X2Fuheun05_126451192236476Sc_cpgX2F250DX250CF258DXX25C4X25BDX25FA_X25C4X25DAX25B7Z25CEX25B325AA_huheun05_1.pjg&type=b400

https://search.static.net/common/?src=http%3A%2F%2Fimages.naver.net%2Fimage%2F081X2F2020X2F01X2F30X2F0003062054_001_20200130114108483_pjg&type=b400

https://search.static.net/common/?src=http%3A%2F%2Fimages.naver.net%2Fimage%2F584X2F2020X2F03X2F05X2F0000007903_001_20200305162019868_pjg&type=b400

https://search.static.net/common/?src=http%3A%2F%2Fimages.naver.net%2Fimage%2F547X2F2020X2F02X2F27X2F0000046762_001_20200227000626470_1.jp&type=b400

https://search.static.net/common/?src=http%3A%2F%2Fimages.naver.net%2Fimage%2F502X2F2020X2F02X2F27X2F0000073802_001_20200221140394900_1.jp&type=b400

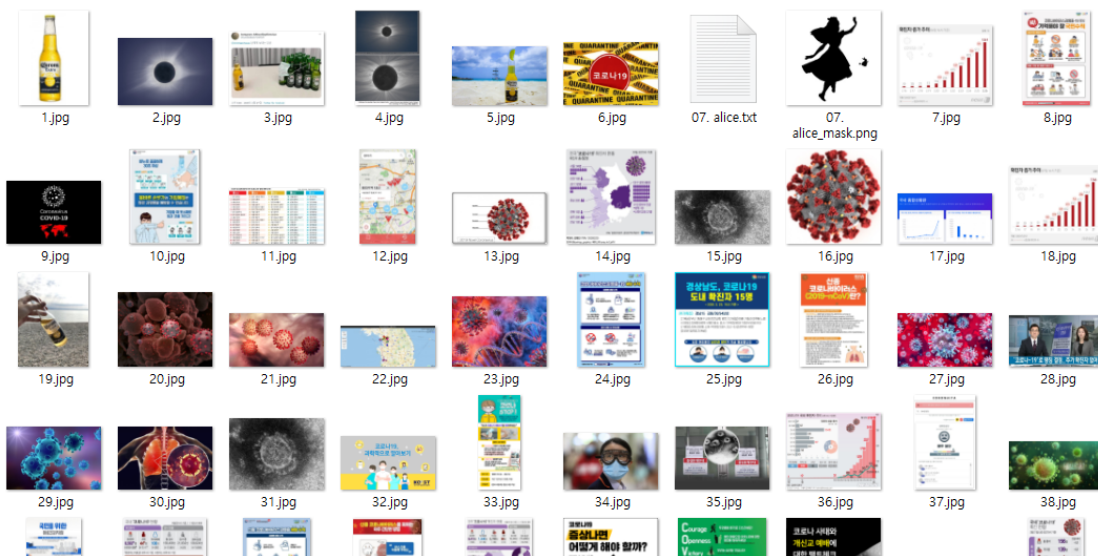
https://search.static.net/common/?src=http%3A%2F%2Fimages.naver.net%2Fimage%2F5461X2F2020X2F04X2F15X2F0000001789_001_20200415091482807_1.jp&type=b400

https://search.static.net/images/search%3A%2F%2Fimages.naver.net%2Fimage%2FD50TTCF0000002FWX2F2WVF00000451_01_20200120092277650_pjg&type=b400

```
In [13]: cnt = 1
for i in img:
    imgUrl = i['data-source']
    # print(imgUrl)
    with urlopen(imgUrl) as file: # 해당 이미지 파일을 저장하기 위한 토끼이다.
        with open('./image/'+str(cnt)+'_jpg','wb') as imgFile: # 해당 파일을 해당 경로에 이름을 만들어 저장한다.wb는 write binary 라는 뜻O.
            img = file.read() # 이미지 파일을 읽음
            imgFile.write(img) # 읽은 파일을 write로 저장한다.
    cnt+=1 # 각 이름을 바꾸기 위해 cnt+1 을 한다.
    if cnt>imgCnt: # 우리가 처음 지정했던 50개가 넘어가면 멈추게 break로 멈추게 한다.
        break
print('Image download success')
```

Image download success

이걸 실행하면 아래와 같이 해당 폴더에 image가 저장되는 것을 볼 수 있다.



'Python' 카테고리의 다른 글

[Python] python 에서 Selenium을 설치 방법

[Python] BeautifulSoup을 통한 이미지 블로그 스크래핑하기

[Python] BeautifulSoup을 통한 이미지 스크래핑 하기

[python] 영화 리뷰에 대한 자연어 처리분석/ 감성분석하기 feat. 스크래핑

[python] BeautifulSoup를 통한 영화리뷰 scraping 하기

[Python] 파이썬 기초 14 - 아주 기초적인 pandas 사용법과 예제

BeutifulSoup 이미지 가져오기

python 이미지 가져오기

python 이미지 스크래핑

python 이미지 크롤링



나무늘보스

혼자 끄적끄적하는 블로그 입니다.