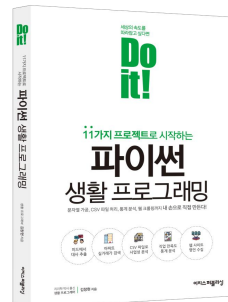


세상의 속도를
따라잡고 싶다면

**Do
it!**

파이썬 생활 프로그래밍



이지스퍼블리싱(주)

06

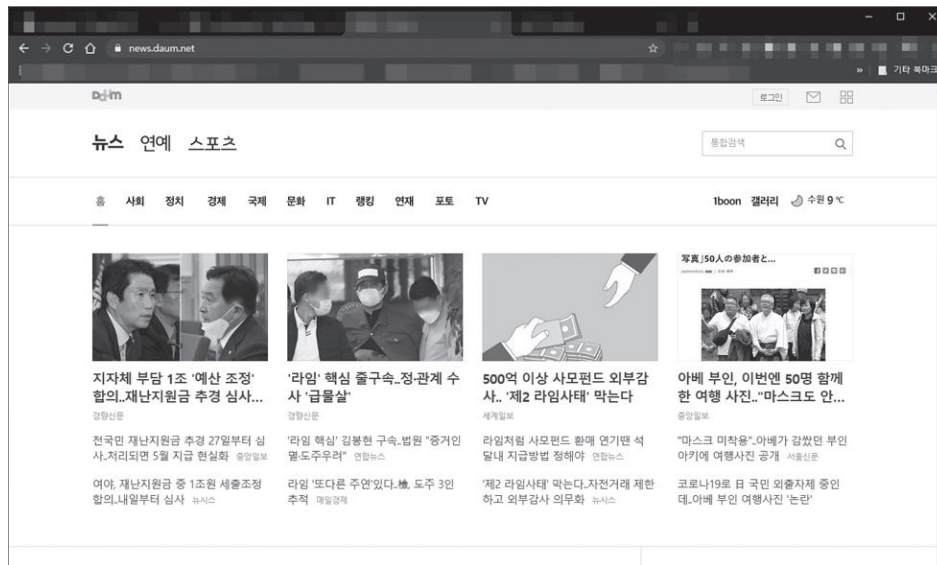
Web Crawling으로 정보 모으기

06-3 Portal Site에서 기사 crawling하기

06-3 Portal Site에서 기사 crawling 하기

이런 상황이라면?

- Portal Site에서 머리기사 정보 모으기



06-3 Portal Site에서 기사 crawling 하기

Web Crawling 기본 환경 준비하기

- 필요한 모듈 임포트

```
>>> import os, re  
>>> import urllib.request as ur  
>>> from bs4 import BeautifulSoup as bs
```

06-3 Portal Site에서 기사 crawling 하기

Web Crawling 기본 환경 준비하기

- 저장할 파일 위치와 접속할 url 저장

```
>>> os.chdir(r'C:\Users\user\python')
```

```
>>> news = 'https://news.daum.net/'
```

06-3 Portal Site에서 기사 crawling 하기

Web Crawling 기본 환경 준비하기

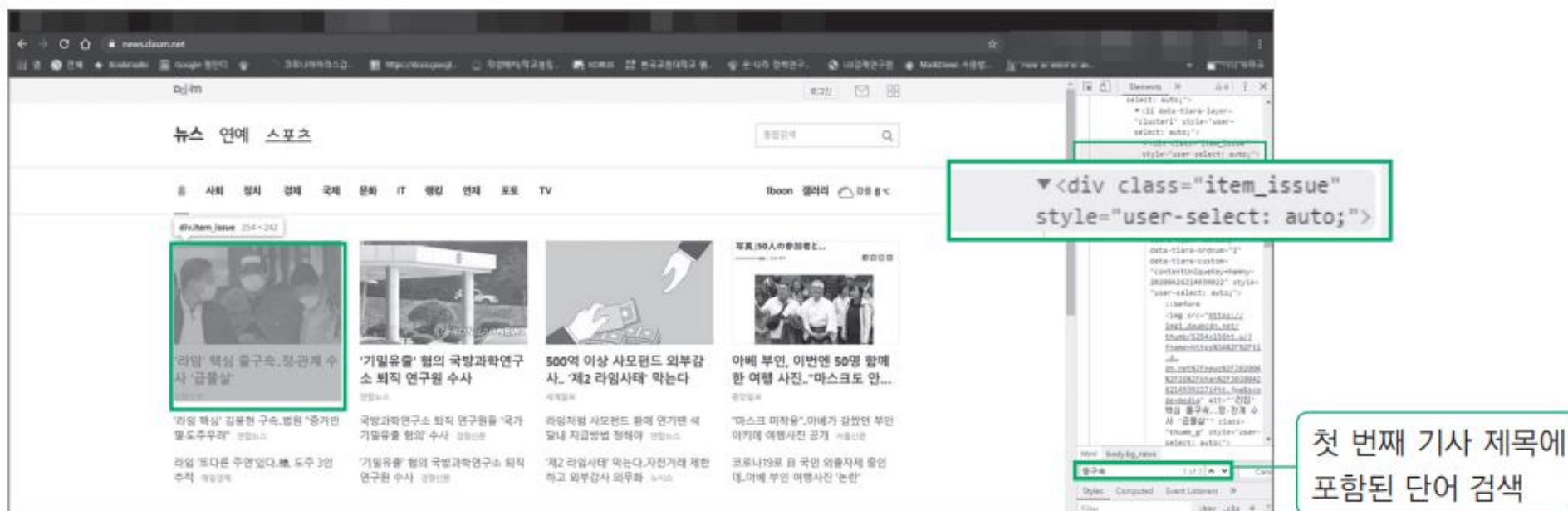
- html 파싱

```
>>> soup = bs(ur.urlopen(news).read(), 'html.parser')
```

06-3 Portal Site에서 기사 crawling 하기

머리기사 제목 추출하기

- 기사 제목에 들어있는 단어 검색



06-3 Portal Site에서 기사 crawling 하기

머리기사 제목 추출하기

- find_all로 <div> 내용 추출

```
soup.find_all('div', {"class" : "item_issue"})
```


06-3 Portal Site에서 기사 crawling 하기

반복문으로 기사 제목 모두 추출하기

- 추출한 div 안의 텍스트를 모두 출력

```
>>> for i in soup.find_all('div', {"class": "item_issue"}):  
    i.text
```

```
'\n\'라임 핵심\' 김봉현 구속..법원 "증거인멸·도주우려"\n연합뉴스\n'
```

```
"\n라임 '또다른 주연'있다..檢, 도주 3인 추적\n매일경제\n"
```

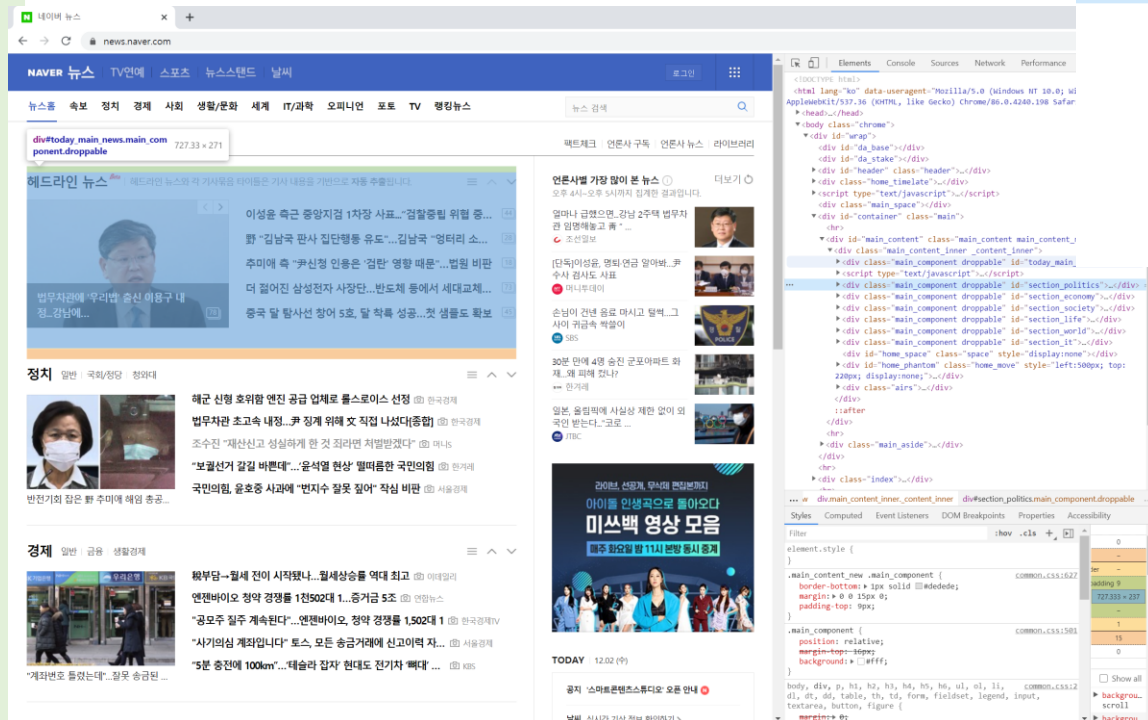
```
"\n국방과학연구소 퇴직 연구원들 '국가 기밀유출 혐의' 수사\n경향신문\n"
```

```
"\n'기밀유출' 혐의 국방과학연구소 퇴직 연구원 수사\n경향신문\n"
```

06-3 Portal Site에서 기사 crawling 하기

머리기사 제목 추출하기

- 기사 제목에 들어있는 단어 검색



```
<div id="container" class="main">
  <hr>
  <div id="main_content" class="main_content main_content_new">
    <div class="main_content_inner _content_inner">
      <div class="main_component droppable" id="today_main_news">...</div>
      <script type="text/javascript">...</script>
      ...
      <div class="main_component droppable" id="section_politics">...</div>
      <div class="main_component droppable" id="section_economy">...</div>
      <div class="main_component droppable" id="section_society">...</div>
      <div class="main_component droppable" id="section_life">...</div>
      <div class="main_component droppable" id="section_world">...</div>
      <div class="main_component droppable" id="section_it">...</div>
      <div id="home_space" class="space" style="display:none">...</div>
      <div id="home_phantom" class="home_move" style="left:500px; top:
        220px; display:none">...</div>
      <div class="airs">...</div>
    </div>
  </div>
</div>
```

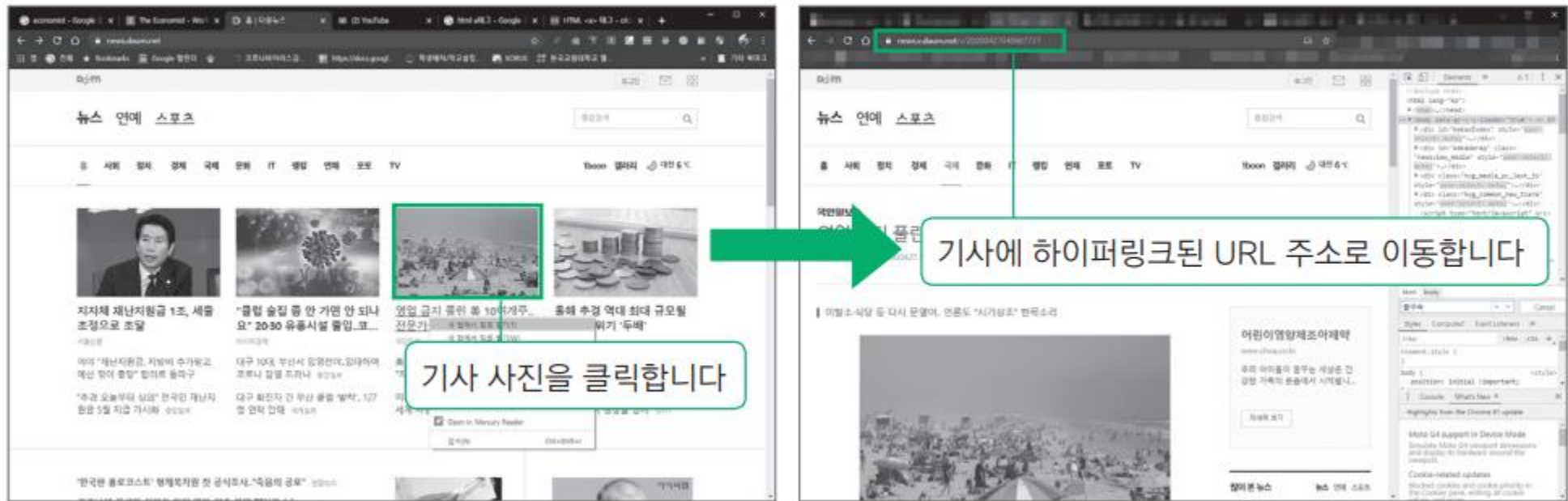
soup.find_all('div', {"class" : "main_component droppable"})
→

```
>>> Warning : ('Connection aborted.', RemoteDisconnected('Remote end closed connection with
  out response'))
>>>
```

06-3 Portal Site에서 기사 crawling 하기

hyperlink 주소 추출하기

- 어떤 요소를 누를 경우 다른 링크로 이동시키는 것이 hyperlink



06-3 Portal Site에서 기사 crawling 하기

hyperlink 주소 추출하기

- <a> tag 사용법 알아보기

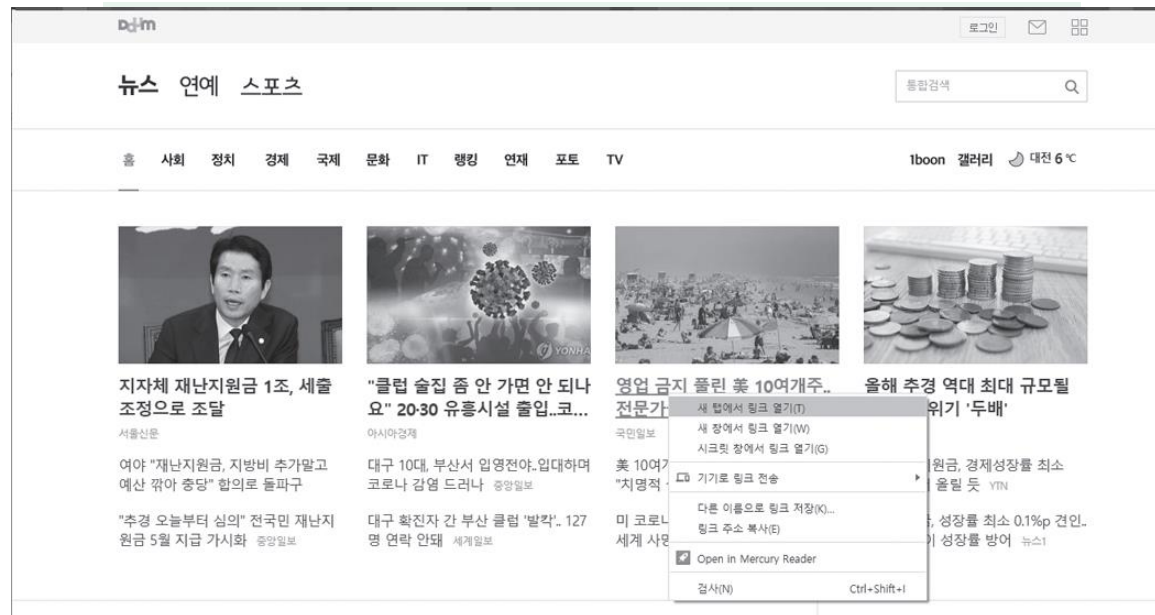
```
<a href = "링크할 URL 주소">하이퍼링크 텍스트</a>
```

```
<HTML>  
  <body>  
    <a href="http://www.naver.com"> 네이버 바로가기 </a>  
  </body>  
</HTML>
```

06-3 Portal Site에서 기사 crawling 하기

hyperlink 주소 추출하기

- 새 탭에서 링크 열기가 있으면 <a> tag를 사용한 것



06-3 Portal Site에서 기사 crawling 하기

hyperlink 주소 추출하기

- <a> tag만 추출하기

```
soup.find_all('a')
```

06-3 Portal Site에서 기사 crawling 하기

hyperlink 주소 추출하기

- href 속성값 추출하기
- 이동할 대상 url을 지정하는 속성
- get() 함수를 통하여 속성 얻기 가능

```
a.get('속성')
```

06-3 Portal Site에서 기사 crawling 하기

hyperlink 주소 추출하기

- href 속성값 추출하기
- 이동할 대상 url을 지정하는 속성
- get() 함수를 통하여 속성 얻기 가능

```
>>> for i in soup.find_all('a')[:5]:  
        i.get('href')
```

```
'#kakaoBody'
```

```
'#kakaoGnb'
```

```
'https://news.daum.net/'
```

```
'https://entertaimdaum.net'
```

```
'https://sports.media.daum.net/sports'
```


06-3 Portal Site에서 기사 crawling 하기

원하는 영역에서 hyperlink 모두 추출하기

- 파싱 가능하게 준비

```
>>> news = 'https://news.daum.net/'  
>>> soup = bs(ur.urlopen(news).read(), 'html.parser')
```

06-3 Portal Site에서 기사 crawling 하기

원하는 영역에서 hyperlink 모두 추출하기

- find_all로 <a> tag 추출하기

```
>>> for i in soup.find_all('div', {"class": "item_issue"}):  
    i.find_all('a')
```

```
[<a class="link_thumb" data-tiara-custom="contentUniqueKey=hamny-20200427060512538"  
data-tiara-id="20200427060512538" data-tiara-layer="article_thumb" data-tiara-ord-  
num="1" data-tiara-type="harmony" href="https://news.v.daum.net/v/20200427060512538">  
  
</a>], <a class="link_txt" data-tiara-custom="contentUniqueKey=hamny-20200427060512538"  
data-tiara-id="20200427060512538" data-tiara-layer="article_main" data-tiara-ord-  
num="1" data-tiara-type="harmony" href="https://news.v.daum.net/v/20200427060512538">  
정부, 무급휴직 주한미군 근로자 임금 선지급 방침</a>]  
[<a class="link_thumb" data-tiara-custom="contentUniqueKey=hamny-20200427060438523"  
data-tiara-id="20200427060438523" data-tiara-layer="article_thumb" data-tiara-ord-  
num="1" data-tiara-type="harmony" href="https://news.v.daum.net/v/20200427060438523">  
...(생략)...
```

06-3 Portal Site에서 기사 crawling 하기

원하는 영역에서 hyperlink 모두 추출하기

- find_all로 <a> tag 추출 후 바로 get은 list 형식이기에 불가능

```
>>> for i in soup.find_all('div', {"class": "item_issue"}):
    i.find_all('a').get('href')

Traceback (most recent call last):
  File "<pyshell#85>", line 2, in <module>
    i.find_all('a').get('href')
  File "C:\Users\user\Anaconda3\lib\site-packages\bs4\element.py", line 1884, in __
    getattr__
    "ResultSet object has no attribute '%s'. You're probably treating a list of items
    like a single item. Did you call find_all() when you meant to call find()?" % key
AttributeError: ResultSet object has no attribute 'get'. You're probably treating a
list of items like a single item. Did you call find_all() when you meant to call find()?
```

06-3 Portal Site에서 기사 crawling 하기

원하는 영역에서 hyperlink 모두 추출하기

- 인덱스를 지정 후 get 사용 가능

```
>>> for i in soup.find_all('div', {"class": "item_issue"}):  
    i.find_all('a')[0]  
  
<a class="link_thumb" data-tiara-custom="contentUniqueKey=hamny-20200427060512538"  
data-tiara-id="20200427060512538" data-tiara-layer="article_thumb" data-tiara-ord-  
num="1" data-tiara-type="harmony" href="https://news.v.daum.net/v/20200427060512538">  
  
</a>  
...(생략)...
```

06-3 Portal Site에서 기사 crawling 하기

원하는 영역에서 hyperlink 모두 추출하기

- 인덱스를 지정 후 get 사용 가능

```
>>> for i in soup.find_all('div',{'class':"item_issue"}):  
    i.find_all('a')[0].get('href')
```

```
'https://news.v.daum.net/v/20200426214939022'
```

```
'https://news.v.daum.net/v/20200426205218584'
```

```
'https://news.v.daum.net/v/20200426201134187'
```

```
'https://news.v.daum.net/v/20200426180251978'
```

06-3 Portal Site에서 기사 crawling 하기

원하는 영역에서 hyperlink 모두 추출하기

- hyperlink를 추출하는 과정

① soup 객체에서 class 속성값이 'item_issue'인 <div> 태그를 find_all로 가져옵니다



② <div> 태그 안에서 <a> 태그를 find_all로 가져옵니다



③ <a> 태그의 href 속성값을 get으로 출력합니다

06-3 Portal Site에서 기사 crawling 하기

기사 제목과 내용 한번에 추출하기

- 기사 URL 얻기



06-3 Portal Site에서 기사 crawling 하기

기사 제목과 내용 한번에 추출하기

- URL 저장 후 BeautifulSoup로 열기

```
>>> article1 = 'https://go.seoul.co.kr/news/newsView.php?id=20200427004004&wlog_tag3=daum'
```

```
>>> soup2 = bs(ur.urlopen(article1).read(), 'html.parser')
```


06-3 Portal Site에서 기사 crawling 하기

기사 제목과 내용 한번에 추출하기

- 기사 내용 가져오기
- <p> tag 안만 확인하면 확인 가능

```
>>> for i in soup2.find_all('p'):  
    print(i.text)
```

오늘 여야 예결특위서 2차 추경 심사...전 국민 지원 확대로 4조 6000억 증액

여야가 27일부터 국회 상임위원회와 예산결산특별위원회를 열고 코로나19 대응 긴급재난지원금 지급을 위한 2차 추가경정예산(추경)안 심사에 들어가기로 했다. 또 전 국민에게 재난지원금을 지급하기 위해 추가로 필요한 자원 가운데 지방정부가 부담할 예정이었던 1조원을 세출 조정을 통해 조달하기로 했다. 더불어민주당 이인영, 미래통합당 심재철 원내대표는 26일 국회에서 이 같은 입장을 밝혔다. 이 원내대표는 “심 원내대표가 ‘지방정부가 당초 부담하기로 했던 1조원 규모라도 세출 조정을 통해 마련하면 어떨겠느냐’고 요청했다”면서 “긴급하게 기획재정부 담당자를 불러 그게 가능한지를 상의했고 최종적으로 가능하게 하기로 정리했다”고 말했다.

...(생략)...

06-3 Portal Site에서 기사 crawling 하기

기사 제목과 내용 한번에 추출하기

- 기사 제목 가져오기
- class 속성이 'item_issue'인 div 안에 존재

```
>>> headline = soup.find_all('div', {"class" : "item_issue"})
```

```
>>> print(headline[0].text)
```

지자체 재난지원금 1조, 세출 조정으로 조달

서울신문

06-3 Portal Site에서 기사 crawling 하기

hyperlink된 모든 기사의 제목과 본문 추출하기

- 머리 기사의 제목 추출

```
>>> for i in headline:  
    print(i.text, '\n')
```

지자체 재난지원금 1조, 세출 조정으로 조달
서울신문

"클럽 술집 좀 안 가면 안 되나요" 20·30 유흥시설 출입...코로나19 확산 우려"
아시아경제

영업 금지 풀린 美 10여개주... 전문가들 "치명적 실수"
국민일보

올해 추경 역대 최대 규모될 듯... 금융위기 '두배'
한국일보

06-3 Portal Site에서 기사 crawling 하기

hyperlink된 모든 기사의 제목과 본문 추출하기

- 기사 제목 출력 후 연결된 링크를 통하여 기사 본문 URL 정보 얻기

```
>>> for i in headline:

    print(i.text, '\n')

    soup3 = bs(ur.urlopen(i.find_all('a')[0].get('href')).read(), 'html.parser')
    for j in soup3.find_all('p'):
        j.text
```

해당 기사의 URL 주소

```

web_crawling_2.py
File Edit Format Run Options Window Help
1 #
2 # 참조 : http://hleecaster.com/python-web-crawling-with-beautifulsoup/
3 # "11가지 프로젝트로 시작하는 생활프로그래밍", 이창현 저, 이지스퍼블리싱, 2020
4 #
5
6 D = True
7 D_1 = False
8
9 import os, re
10 import usecsv
11 import requests
12 import urllib.request as ur
13
14 from bs4 import BeautifulSoup as bs
15
16 try:
17     os.chdir(r'C:\과소사\과소사-강의예제\web_crawling')
18
19     """
20     news = "https://news.yahoo.com/"
21     news = "https://www.chosun.com/"
22     news = "https://news.naver.com/"
23     """
24     news = "https://news.daum.net/"
25
26     webpage = requests.get(news)
27     if D:
28         print("\n1) >> webpage : ", webpage)
29
30     soup = bs(webpage.content, 'html.parser')
31     if D:
32         print("\n2) >> soup : ", soup)
33
34     # 기사 제목 추출하기
35     # find_all로 <div> 내용 추출
36     # class 속성이 'item_issue'인 div 안에 존재
37     if D:
38         print("\n4) >> 기사 제목 추출하기")
39
40     headline = soup.find_all('div', {"class" : "item_issue"})
41
42     for i in headline:
43         print(i.text, "\n")
44
45     # find_all로 <a> tag 추출하기
46     if D:
47         print("\n5) >> find_all로 <a> tag 추출하기 : ")
48
49     for i in soup.find_all('a')[:5]:
50         print(i.get('href'))
51

```

```

52 # 원하는 영역에서 하이퍼링크 모두 추출하기
53 # 인덱스를 지정 후 get 사용 가능
54 if D:
55     print("\n6) >> 원하는 영역에서 하이퍼링크 모두 추출하기")
56
57 for i in headline:
58     print(i.find_all('a')[0].get('href'))
59
60
61 # 기사 제목 출력 후 연결된 링크를 통하여 기사 본문 URL 정보 얻기
62 if D:
63     print("\n7) >> 기사 제목 출력 후 연결된 링크를 통하여 기사 본문 URL 정보 얻기")
64 for i in headline:
65     if D:
66         print("\n7-1) >> 상위기사 제목 : ")
67         print(i.text, "\n")
68
69         new_link = i.find_all('a')[0].get('href')
70         link_webpage = requests.get(new_link)
71         if D:
72             print("\n7-2) >> new_link : ", new_link)
73             print("\n7-3) >> link_webpage : ", link_webpage)
74
75         soup_link = bs(link_webpage.content, 'html.parser')
76         for j in soup_link.find_all('p'):
77             print(j.text, "\n")
78
79 # Portal Site에서 기사 crawling 하기
80 # 기사 제목과 내용 한번에 추출하기
81 if D:
82     print("\n8) >> Portal Site에서 기사내용 crawling 하기")
83     print("\n      기사 제목과 내용 한번에 추출하기")
84
85     article_link = 'http://go.seoul.co.kr/news/newsView.php?id=20201109014005'
86
87     article = requests.get(article_link)
88     if D:
89         print("\n9) >> article : ", article)
90
91     soup_article = bs(article.content, 'html.parser')
92     if D:
93         print("\n10) >> soup_article : ", soup_article)
94         print("\n11) >> 전체기사 내용")
95
96     for i in soup_article.find_all('p'):
97         print(i.text)
98
99
100 except Exception as e:
101     print("\n>>> Warning : ", e)
102

```

```
Python 3.8.5 Shell
File Edit Shell Debug Options Window Help
Python 3.8.5 (tags/v3.8.5:580fbb0, Jul 20 2020, 15:43:08) [MSC v.1926 32 bit (Intel)] on win32
Type "help", "copyright", "credits" or "license()" for more information.
>>>|
===== RESTART: C:\W과소사\W과소사-강의예제\Wweb_crawling\Wweb_crawling_2.py ==
=====

1) >> webpage : <Response [200]>

2) >> soup : Squeezed text (1204 lines).

4) >> 기사 제목 추출하기

문대통령 "코로나 극복의지 담은 예산 통과.. 여야에 감사"

국민일보

복귀한 윤석열, 월성 원전 영장부터 승인..파장

머니투데이

조슈아 원 13.5개월..'홍콩 민주화' 청년 3명 실행 선고

한겨레
```

```
Python 3.8.5 Shell
File Edit Shell Debug Options Window Help
로써 금연 정책에 대한 변혁이 시작될 것으로 기대하고 있다"면서 "앞으로 방배동, 서초동, 반포동, 잠원동 등으로 확대하겠다"고 밝혔다.이민영 기자 min@seoul.co.kr

2020-11-09 14면

김인호 서울시의회 의장, 온라인 청소년
김인호 서울시의회 의장, 안전한 수능 실
서울특별시의회 김인호 의장(더불어민주당, 동대문3)은 오는 12월 3일에 치러지는 2021학년도 대학수
학능력시험의...
정윤경 경기도의회 교육위원장,도교육청
장현국 경기도의회 의장, 일일 소원으로

조은희 서초구청장 서울시장 출마 "여성가산점 안 받고 실력
"지금은 남성·여성보다 일 잘하는 일꾼 필요"
정무부시장·구청장 등 서울행정 10년 경험
내일 부동산·세금 문제 등 입장 발표 예정
김종인 "文정부 비판보다 시민 마음 얻길"

"공공원룸 배란다는 주거인권... 국유지에 주택 공급"
쪽방촌 재개발하는 김영종 종로구청장

수험생 지원!... 광진, 고3 1인당 마스크 10장씩
학원·교습소 등 815곳도 16만장 전달
수능 당일 수험생 수송 상황실 운영

"장애인 배려·주민 편의 원원 복지관"
[현장 행정] 은평 2호 '우리장애인복지관' 개관

최신 장비 시설로 장애인들 복지 향상
주민 편의시설 체력단련실·카페 갖춰
초기 주민들 반대 어려움 딛고 문열어
김미경 구청장 "장애인 행복한 삶 기여"

은평, 문체부 등 평가·공모사업 성적 탁월
금천구의 자랑, 청소년상담복지센터 '안전
'라이브 관악' 구독·인증 땀 추첨 통해
코로나 블루 날리는 강동 '희망의 빛'...
수험생 지원!... 광진, 고3 1인당 마스크 10
중구 주민 여러분, 마음 안녕하신가요

자료 제공 : 정책브리핑 korea.kr
주소 : 100-745 서울시 중구 세종대로 124 (태평로1가 25번지) 서울신문사빌딩 | 대표전화 : (02) 20
00-9000
인터넷서울신문에 게재된 콘텐츠의 무단 전재/복사/배포 행위는 저작권법에 저촉되며 위반 시 법적 제
재를 받을 수 있습니다.
Copyright © 서울신문사 All rights reserved.
>>>
```

감사합니다

