

세상의 속도를
따라잡고 싶다면

**Do
it!**

파이썬 생활 프로그래밍



이지스퍼블리싱(주)

06

Web Crawling으로 정보 모으기

06-4 Program 실행 file 만들기

06-4 Program 실행 file 만들기

URL 주소 저장하기

- 출력한 기사 URL 주소 저장하기

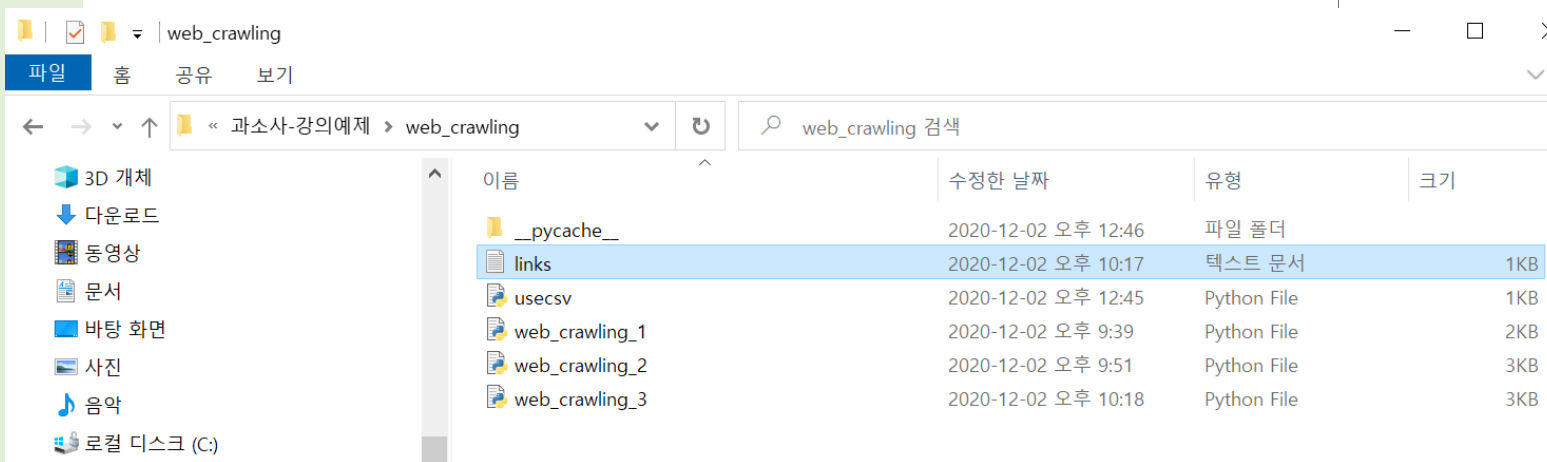
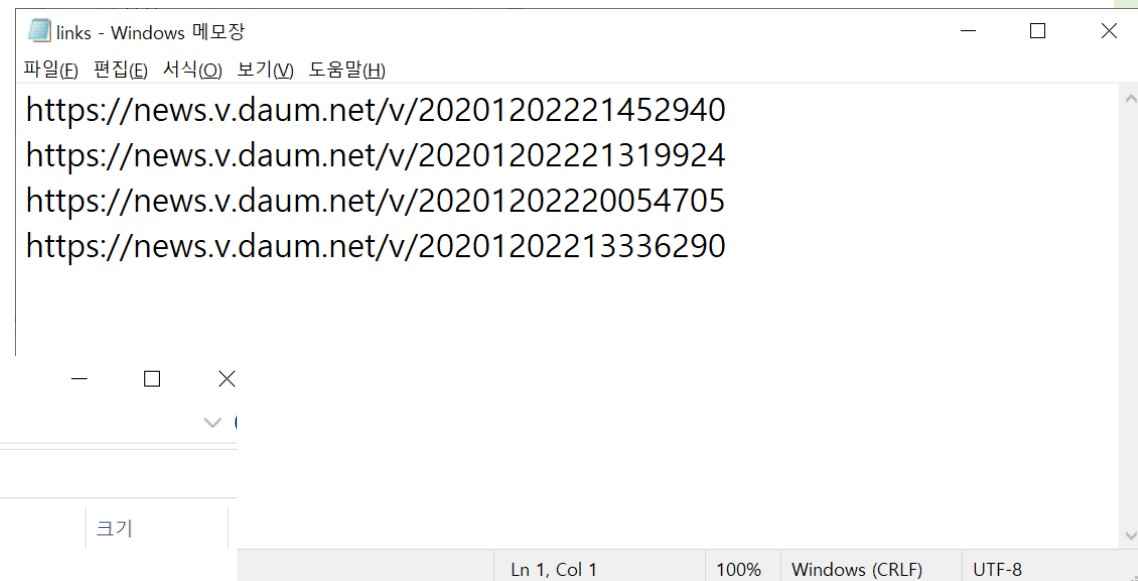
```
>>> os.chdir(r'C:\Users\user\python')
>>> f = open('links.txt', 'w')
>>> for i in soup.find_all('div', {"class": "item_issue"}):
    f.write(i.find_all('a')[0].get('href')+'\n' )

44
44
44
44
>>> f.close()
```

06-4 Program 실행 file 만들기

URL 주소 저장하기

- 저장한 file 확인



06-4 Program 실행 file 만들기

기사 본문을 file로 저장하기

- Python으로 기사의 내용 저장

```
>>> article1 = 'https://news.v.daum.net/v/20200430135751773'

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

>>> f = open('article_1.txt', 'w')

>>> for i in soup.find_all('p'):
        f.write(i.text)

85
100
21
0
203
249
109
177
183
112
26
>>> f.close()
```

06-4 Program 실행 file 만들기

기사 본문을 file로 저장하기

- 저장한 file 확인
- → 다음 페이지 예제와 통합해서 실행

06-4 Program 실행 file 만들기

기사 제목, 본문, hyperlink를 file로 저장하기

- 제목, 링크, 내용 순으로 저장
- File open 시 'w' 또는 'a' mode

```
>>> url = 'https://news.daum.net/'
>>> soup = bs(ur.urlopen(url).read(), 'html.parser')
>>> f = open('article_total.txt', 'w')
>>> for i in soup.find_all('div', {"class": "item_issue"}):
    try:

        f.write(i.text + '\n')

        f.write(i.find_all('a')[0].get('href') + '\n')

        soup2 = bs(ur.urlopen(i.find_all('a')[0].get('href')).read(), 'html.parser')

        for j in soup2.find_all('p'):
            f.write(j.text)

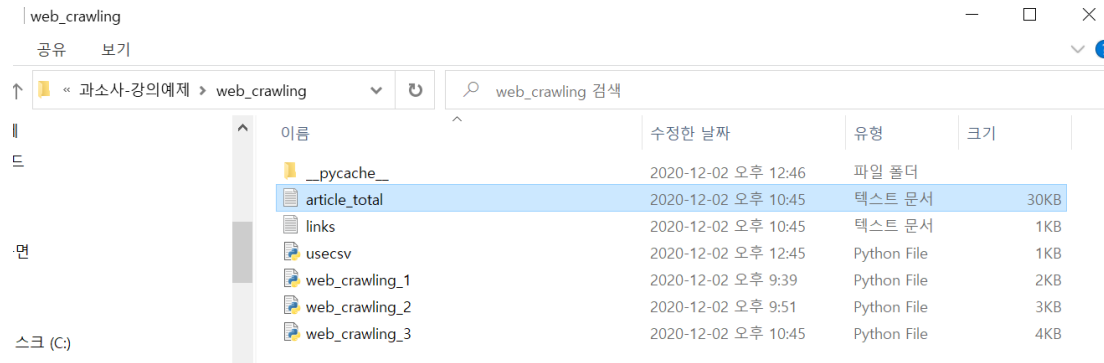
    except:
        pass

>>> f.close()
```

06-4 Program 실행 file 만들기

기사 제목, 본문, hyperlink를 file로 저장하기

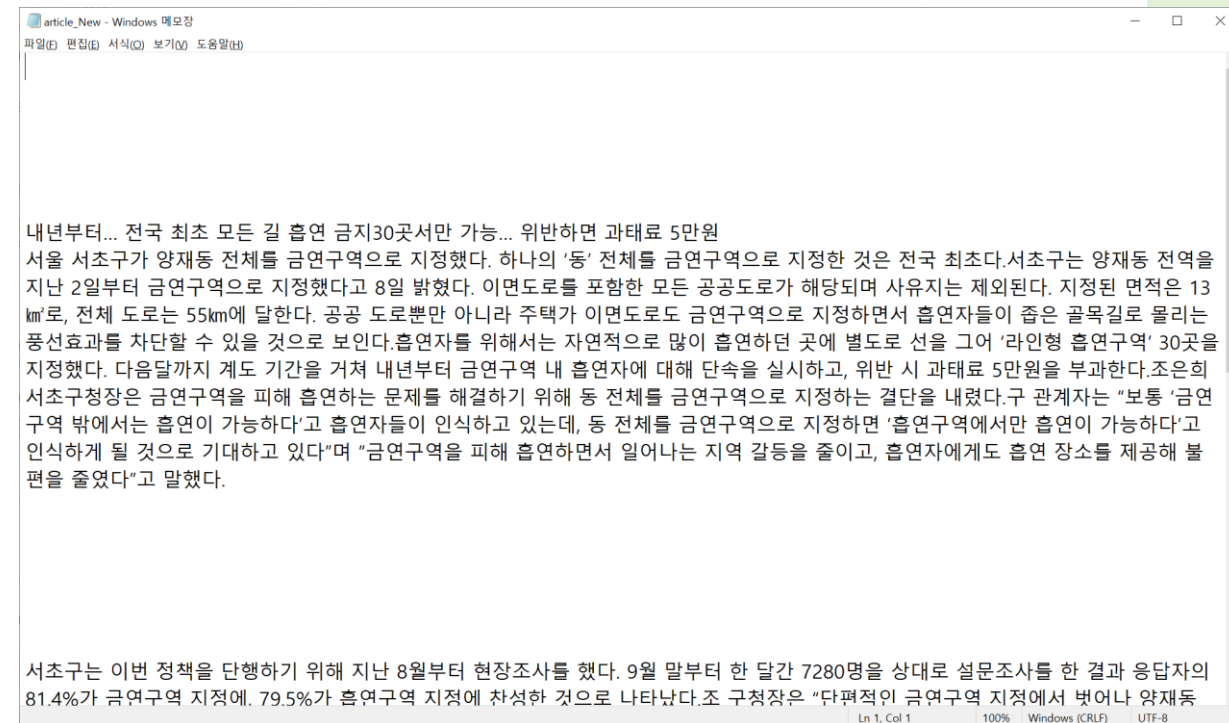
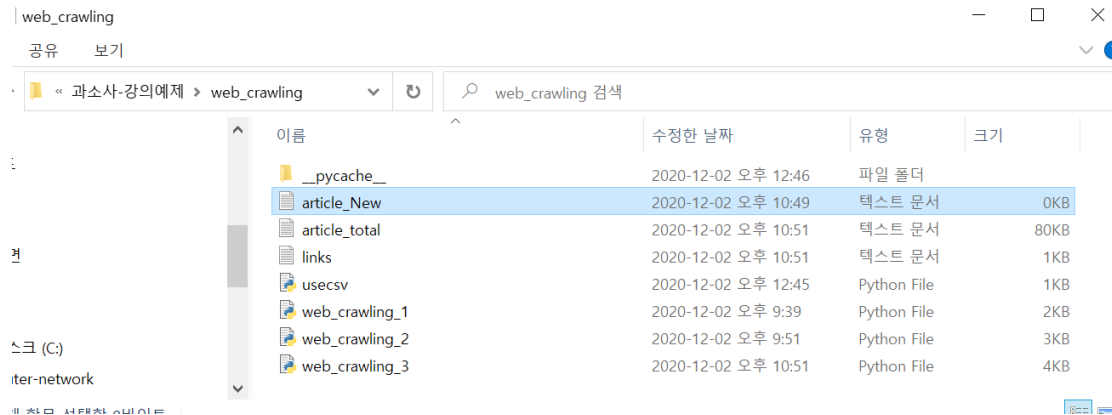
- 저장한 file 확인



06-4 Program 실행 file 만들기

기사 제목, 본문, hyperlink를 file로 저장하기

- 저장한 file 확인



```

web_crawling_3.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     news = "https://news.daum.net/"
20
21     webpage = requests.get(news)
22     if D:
23         print("\n1) >> webpage : ", webpage)
24
25     soup = bs(webpage.content, 'html.parser')
26     if D:
27         print("\n2) >> soup : ", soup)
28
29     # 기사 제목 추출하기
30     # find_all로 <div> 내용 추출
31     # class 속성이 'item_issue'인 div 안에 존재
32     if D:
33         print("\n4) >> 기사 제목 추출하기")
34
35     headline = soup.find_all('div', {"class" : "item_issue"})
36
37     for i in headline:
38         print(i.text, "\n")
39
40     # find_all로 <a> tag 추출하기
41     if D:
42         print("\n5) >> find_all로 <a> tag 추출하기 : ")
43
44     for i in soup.find_all('a')[:5]:
45         print(i.get('href'))
46
47     # 원하는 영역에서 하이퍼링크 모두 추출하기
48     # 인덱스를 지정 후 get 사용 가능
49     # 추출한 hyperlink들을 file f에 저장
50     if D:
51         print("\n6) >> 원하는 영역에서 하이퍼링크 모두 추출하여 file로 저장")
52

```

```

53 # 한글처리시 error 발생가능. UNICODE로 처리
54 # -1 은 buffer
55 f = open('links.txt', 'w', -1, "utf-8")
56
57 for i in headline:
58     print(i.find_all('a')[0].get('href'))
59     f.write(i.find_all('a')[0].get('href')+'\n')
60
61 f.close()
62
63 # 기사 제목 출력 후 연결된 링크를 통하여 기사 본문 URL 정보 얻은 후 file에 저장
64 if D:
65     print("\n7) >> 기사 제목 출력 후 연결된 링크를 통하여 기사 본문 URL 정보 얻은 후 W
66 file에 저장")
67
68 f_1 = open('article_total.txt', 'a', -1, "utf-8")
69
70 for i in headline:
71     if D:
72         print("\n7-1) >> 상위기사 제목 : ")
73         print(i.text, "\n")
74         f_1.write(i.text)
75
76     new_link = i.find_all('a')[0].get('href')
77     link_webpage = requests.get(new_link)
78     if D:
79         print("\n7-2) >> new_link : ", new_link)
80         print("\n7-3) >> link_webpage : ", link_webpage)
81
82     soup_link = bs(link_webpage.content, 'html.parser')
83     for j in soup_link.find_all('p'):
84         print(j.text, "\n")
85         f_1.write(j.text)
86
87 f_1.close()
88

```

```

89 # Portal Site에서 기사 crawling 하기
90 # 기사 제목과 내용 한번에 추출하여 file에 저장
91 if D:
92     print("\n8) >> Portal Site에서 기사내용 crawling 하기")
93     print("\n      기사 제목과 내용 한번에 추출하여 file에 저장 ")
94
95 article_link = 'http://go.seoul.co.kr/news/newsView.php?id=20201109014005'
96
97 article = requests.get(article_link)
98 if D:
99     print("\n9) >> article : ", article)
100
101 soup_article = bs(article.content, 'html.parser')
102 if D:
103     print("\n10) >> soup_article : ", soup_article)
104     print("\n11) >> 전체기사 내용")
105
106 f_2 = open('article_New.txt', 'w', -1, "utf-8")
107
108 for i in soup_article.find_all('p'):
109     print(i.text)
110     f_2.write(i.text)
111
112 f_2.close()
113
114
115 except Exception as e:
116     print("\n>>> Warning : ", e)

```

```

Python 3.8.5 Shell
File Edit Shell Debug Options Window Help

학능력시험의...
정운경 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”

“공공원을 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 자량, 청소년상담복지센터 ‘안전
‘라이브 관악’ 구독·인증 땀 추첨 통해
코로나 블루 날리는 강동 ‘희망의 빛’...
수험생 지원!… 광진, 고3 1인당 마스크 10
중구 구민 여러분, 마음 안녕하십니까

자료 제공 : 정책브리핑 korea.kr
주소 : 100-745 서울시 중구 세종대로 124 (태평로1가 25번지) 서울신문사빌딩 I 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.
>>>

```

Ln: 2639 Col: 4

06-4 Program 실행 file 만들기

Web Crawling 실행 file 만들기

- pip로 pyinstaller 설치하기

```
C:\Users\user>pip install pyinstaller
```

```
Windows PowerShell
PS C:\Users\user>
PS C:\Users\user> pip install pyinstaller
Collecting pyinstaller
  Downloading pyinstaller-4.1.tar.gz (3.5 MB)
    | 3.5 MB 1.7 MB/s
  Installing build dependencies ... done
  Getting requirements to build wheel ... done
  Preparing wheel metadata ... done
Requirement already satisfied: setuptools in c:\users\user\AppData\Local\programs\python\python38-32\lib\site-packages (from pyinstaller) (47.1.0)
Collecting pywin32-ctypes>=0.2.0; sys_platform == "win32"
  Downloading pywin32-ctypes-0.2.0-py2.py3-none-any.whl (28 kB)
Collecting altgraph
  Downloading altgraph-0.17-py2.py3-none-any.whl (21 kB)
Collecting pyinstaller-hooks-contrib>=2020.6
  Downloading pyinstaller_hooks_contrib-2020.10-py2.py3-none-any.whl (166 kB)
    | 166 kB ...
Collecting pefile>=2017.8.1; sys_platform == "win32"
  Downloading pefile-2019.4.18.tar.gz (62 kB)
    | 62 kB 139 kB/s
Collecting future
  Downloading future-0.18.2.tar.gz (829 kB)
    | 829 kB 6.8 MB/s
Using legacy 'setup.py install' for pefile, since package 'wheel' is not installed.
Using legacy 'setup.py install' for future, since package 'wheel' is not installed.
Building wheels for collected packages: pyinstaller
  Building wheel for pyinstaller (PEP 517) ... done
  Created wheel for pyinstaller: filename=pyinstaller-4.1-py3-none-any.whl size=2790249 sha256=fbf803fa0aaad81968e10b2027392cb497a245def8103dfb0cbc7596f086f1c9
  Stored in directory: c:\users\user\AppData\Local\pip\cache\wheels\ae\7a\1e\42202ec16f036e6c25592c6bc63d3c26e6a6a
  ddd6a25f053a
Successfully built pyinstaller
Installing collected packages: pywin32-ctypes, altgraph, pyinstaller-hooks-contrib, future, pefile, pyinstaller
  Running setup.py install for future ... done
  Running setup.py install for pefile ... done
Successfully installed altgraph-0.17 future-0.18.2 pefile-2019.4.18 pyinstaller-4.1 pyinstaller-hooks-contrib-2020.10 pywin32-ctypes-0.2.0
WARNING: You are using pip version 20.2.4; however, version 20.3 is available.
You should consider upgrading via the 'c:\users\user\AppData\Local\programs\python\python38-32\python.exe -m pip install --upgrade pip' command.
PS C:\Users\user>
```

06-4 Program 실행 file 만들기

Web Crawling 실행 file 만들기

- 작성한 코드 Python(.py) file로 저장

```
import os , codecs, re, datetime, requests
import urllib.request as ur
from bs4 import BeautifulSoup as bs
os.chdir(r'C:\Users\user\python')
url = 'https://news.daum.net/'
f = open(str(datetime.date.today()) + 'articles.txt', 'w')
soup = bs(ur.urlopen(url).read(), 'html.parser')

for i in soup.find_all('div', {"class": "thumb_relate"}):
    try:
        f.write(i.text + '\n')
        f.write(i.find_all('a')[0].get('href') + '\n')
        soup2 = bs(ur.urlopen(i.find_all('a')[0].get('href')).read(), 'html.
parser')
        for j in soup2.find_all('p'):
            f.write(j.text)
    except:
        pass

f.close()
```

06-4 Program 실행 file 만들기

Web Crawling 실행 file 만들기

- Python file 실행 file로 만들기

```
pyinstaller --onefile [파이썬 파일].py
```

```
pyinstaller --onefile article_collector.py
```

>> pyinstaller --onefile web_crawling_3.py

```
Windows PowerShell
PS C:\#\과사#\과사-강의예제#\web_crawling> ls

디렉터리: C:\#\과사#\과사-강의예제#\web_crawling

Mode                LastWriteTime         Length Name
----                -
d-----         2020-12-02 오후 12:46             __pycache__
-a-----         2020-12-02 오후 10:59             4320 article_New.txt
-a-----         2020-12-02 오후 10:59          138501 article_total.txt
-a-----         2020-12-02 오후 10:58             180 links.txt
-a-----         2020-12-02 오후 12:45             621 usecsv.py
-a-----         2020-12-02 오후 9:39             1786 web_crawling_1.py
-a-----         2020-12-02 오후 9:51            3047 web_crawling_2.py
-a-----         2020-12-02 오후 10:58            3477 web_crawling_3.py

PS C:\#\과사#\과사-강의예제#\web_crawling> pyinstaller --onefile web_crawling_3.py
62 INFO: PyInstaller: 4.1
62 INFO: Python: 3.8.5
62 INFO: Platform: Windows-10-10.0.18362-SPO
62 INFO: wrote C:\#\과사#\과사-강의예제#\web_crawling#\web_crawling_3.spec
62 INFO: UPX is not available.
62 INFO: Extending PYTHONPATH with paths
['C:\#\과사#\과사-강의예제#\web_crawling', 'C:\#\과사#\과사-강의예제#\web_crawling']
78 INFO: checking Analysis
78 INFO: Building Analysis because Analysis-00.toc is non existent
78 INFO: Initializing module dependency graph...
78 INFO: Caching module graph hooks...
93 INFO: Analyzing base_library.zip ...
2294 INFO: Processing pre-find module path hook distutils from 'c:\#\users#\user#\appdata#\local#\programs#\python#\python38-32#\lib#\site-packages#\PyInstaller#\hooks#\pre_find_module_path#\hook-distutils.py'.
2295 INFO: distutils: retargeting to non-venv dir 'c:\#\users#\user#\appdata#\local#\programs#\python#\python38-32#\lib'
5412 INFO: Caching module dependency graph...
5537 INFO: running Analysis Analysis-00.toc
5537 INFO: Adding Microsoft.Windows.Common-Controls to dependent assemblies of final executable
   required by c:\#\users#\user#\appdata#\local#\programs#\python#\python38-32#\python.exe
5647 INFO: Analyzing C:\#\과사#\과사-강의예제#\web_crawling#\web_crawling_3.py
5771 INFO: Processing pre-safe import module hook urllib3.packages.six.moves from 'c:\#\users#\user#\appdata#\local#\programs#\python#\python38-32#\lib#\site-packages#\PyInstaller#\hooks#\pre_safe_import_module#\hook-urllib3.packages.six.moves.py'.
7219 INFO: Processing module hooks...
```

```
Windows PowerShell
PS C:\#\과사#\과사-강의예제#\web_crawling>
PS C:\#\과사#\과사-강의예제#\web_crawling>
PS C:\#\과사#\과사-강의예제#\web_crawling>
PS C:\#\과사#\과사-강의예제#\web_crawling> ls

디렉터리: C:\#\과사#\과사-강의예제#\web_crawling

Mode                LastWriteTime         Length Name
----                -
d-----         2020-12-02 오후 11:08             build
d-----         2020-12-02 오후 11:08             dist
d-----         2020-12-02 오후 11:09             __pycache__
-a-----         2020-12-02 오후 10:59             4320 article_New.txt
-a-----         2020-12-02 오후 10:59          138501 article_total.txt
-a-----         2020-12-02 오후 10:58             180 links.txt
-a-----         2020-12-02 오후 12:45             621 usecsv.py
-a-----         2020-12-02 오후 9:39             1786 web_crawling_1.py
-a-----         2020-12-02 오후 9:51            3047 web_crawling_2.py
-a-----         2020-12-02 오후 10:58            3477 web_crawling_3.py
-a-----         2020-12-02 오후 11:08             927 web_crawling_3.spec

PS C:\#\과사#\과사-강의예제#\web_crawling> cd dist
PS C:\#\과사#\과사-강의예제#\web_crawling#\dist> ls

디렉터리: C:\#\과사#\과사-강의예제#\web_crawling#\dist

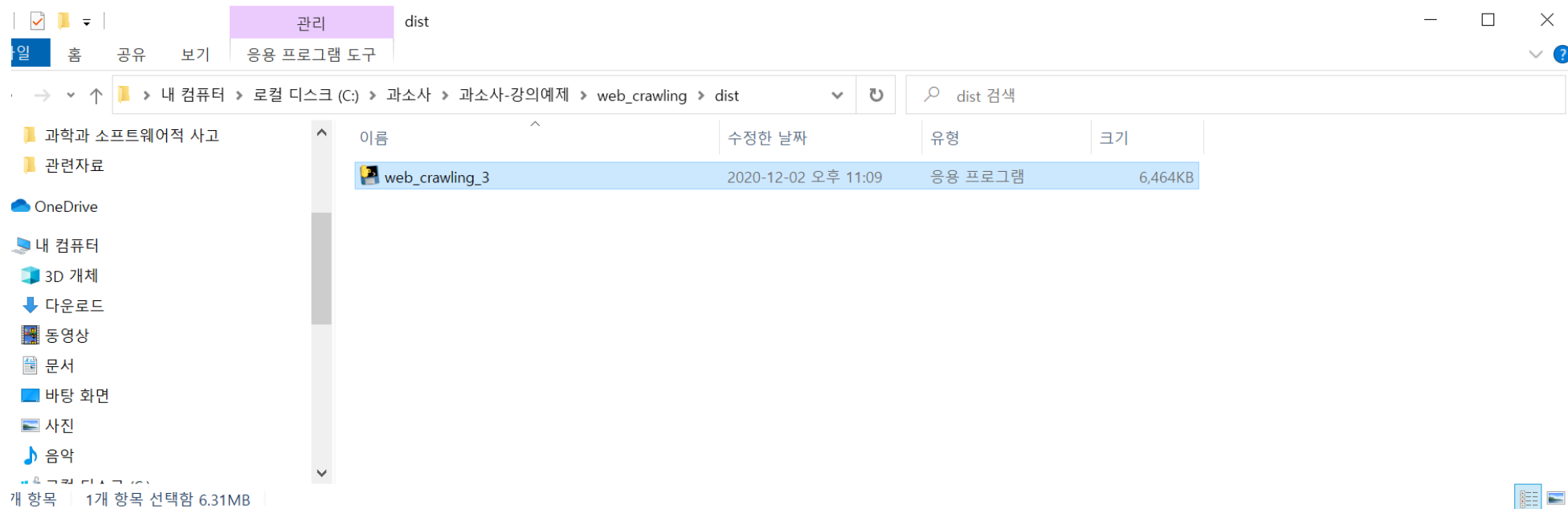
Mode                LastWriteTime         Length Name
----                -
-a-----         2020-12-02 오후 11:09          6618826 web_crawling_3.exe

PS C:\#\과사#\과사-강의예제#\web_crawling#\dist>
```

06-4 Program 실행 file 만들기

Web Crawling 실행 file 만들기

- 'web_crawling_3.exe' 가 생성



06-4 Program 실행 file 만들기

Web Crawling 실행 file 만들기

- web_crawling_3.exe 실행
- Program01 crawling 후 앞에서 실행시 생성하였던 file 생성

감사합니다

