

1-1. 오픈소스 패키지

1-2. 데이터 수집기 완성하기

2. 안티크롤링과 회피법

3. 선택자 심화와 수집기 완성하기

4. url과 요청값 이해하기

1. 오픈소스 패키지

대표적인 파이썬 패키지:

django

Flask

OpenPyXL

Requests

Se

1-2. 데이터 수집기 완성하기

- Requests를 활용한 코드

week3_1.py

```
1 import requests
2 from bs4 import BeautifulSoup
3
4 raw = requests.get("https://tv.naver.com/r")
5 print(raw)
```

1-2. 데이터 수집기 완성하기

- BeautifulSoup4을 활용한 코드

week3_1.py

```
1 import requests
2 from bs4 import BeautifulSoup
3
4 raw = requests.get("https://tv.naver.com/r")
5 # print(raw.text)
6 # 소스코드 출력부분을 주석처리해줍니다!
7
8 html = BeautifulSoup(raw.text, "html.parser")
9 print(html)
```

#1 컨테이너 수집 → #2 영상 별 데이터 수집 → #3 수집 반복

week3_1.py

```
1 import requests
2 from bs4 import BeautifulSoup
3
4 raw = requests.get("https://tv.naver.com/r")
5 html = BeautifulSoup(raw.text, "html.parser")
6
7 clips = html.select("div.inner")
8 # print(clips[0])
9 # 컨테이너 소스코드는 주석처리해줍니다.
10
11 for cl in clips:
12     title = cl.select_one("dt.title")
13     chn = cl.select_one("dd.chn")
14     hit = cl.select_one("span.hit")
15     like = cl.select_one("span.like")
16
17     print(title.text.strip())
18     print(chn.text.strip())
19     print(hit.text.strip())
20     print(like.text.strip())
```

=====

[illegible]

```
1 import requests
2 from bs4 import BeautifulSoup
3
4 raw = requests.get("https://search.naver.com/search.naver?where=news&query=코알라",
5                   headers={'User-Agent': 'Mozilla/5.0'})
6 html = BeautifulSoup(raw.text, "html.parser")
```

#1 컨테이너 수집 → #2 기사 별 기사제목과 언론사 데이터 수집 → #3 수집 반복

```

1 import requests
2 from bs4 import BeautifulSoup
3
4 raw = requests.get("https://search.naver.com/search.naver?where=news&query=코알라",
5                   headers={'User-Agent': 'Mozilla/5.0'})
6 html = BeautifulSoup(raw.text, "html.parser")
7
8 articles = html.select("ul.type01 > li")
9
10 # 리스트를 사용한 반복문으로 모든 기사에 대해서 제목/언론사 출
11 for ar in articles:
12     title = ar.select_one("a._sp_each_title").text
13     source = ar.select_one("span._sp_each_source").text
14
15     print(title, source)

```

4. URL과 요청값 이해하기

URL 요청값의 기본구조

URL 주소는 사람마다 다를 수 있습니다.

홈페이지 주소: 페이지의 기본 주소입니다.

?: 요청값과 홈페이지 주소를 연결해줍니다.

요청값: 요청값에 따라 페이지에 표시되는 내용이 바뀝니다.

where=카테고리

query=검색어

start=기사번호

해당하는 카테고리의 검색결과를 보여줍니다.

해당하는 검색어의 검색결과를 보여줍니다.

해당하는 기사번호의 기사부터 보여줍니다.

where=news

query=코알라

start=11

뉴스 1-10 / 7,686건

&start=1

뉴스 11-20 / 7,689건

&start=11

```
import requests
```

```
from bs4 import BeautifulSoup
```

```
for n in range(1, 100, 10):
```

```
    raw = requests.get("https://search.naver.com/search.naver?where=news&query=코알라  
&start="+str(n), headers={'User-Agent':'Mozilla/5.0'})
```

```
    html = BeautifulSoup(raw.text, "html.parser")
```

```
    articles = html.select("ul.type01 > li"
```

```
    for ar in articles:
```

```
        title = ar.select_one("a._sp_each_title").text
```

```
        source = ar.select_one("span._sp_each_source").text
```

```
        print(title, source)
```