

*2020 D&A*

*DEEP SESSION*

3. 크롤링

# Contents

- 01 크롤링(Crawling)이란?
- 02 HTML의 구조
- 03 Requests와 BeautifulSoup4
- 04 Selenium
- 05 데이터 저장하기

## 01 크롤링(Crawling)이란?

소프트웨어가 웹을 돌아다니며  
필요한 정보를 긁어오는 것을 말함  
(web scraping)

웹페이지는 html의 태그로 구성되어 있음!  
크롤링은 각각의 태그에 저장되어 있는  
데이터를 긁어오는 일을 함

F12를 누르면 확인 가능!

The screenshot shows the Daum homepage with the developer tools (F12) open on the right. The HTML structure is visible, showing the `<body>` tag and its children, including the `daumIndex` and `daumWrap` divs. The `daumWrap` div contains several script tags and an `iframe` for the quicksearch function. The `iframe` has attributes like `frameborder="0"`, `scrolling="no"`, `marginwidth="0"`, `marginheight="0"`, `vspace="0"`, `hspace="0"`, `allowtransparency="true"`, `id="adfit_frame_1c73937df3f"`, and `name="..."`. The `src` attribute is set to `https://t1.daumcdn.net/b2/ssp/awsa.html`. The `style` attribute is `width: 0px; height: 0px; border: 0px; display: none;`. The `whale-quicksearch` function is also visible in the HTML.

02 HTML의 구조    일반적으로 태그는 <html></html> 구조

정보를 찾는 과정  
<h1>D&A</h1>  
-> h1태그에 D&A라는 정보가 들어있음!

태그의 종류    (우리는 이를 이용해 정보를 찾을 것이니 외울 필요는 없음)



태그	설명	예	속성
h1, h2, h3, h4, h5, h6	글자의 크기를 조절	<code>&lt;h1&gt;D&amp;A&lt;/h1&gt;</code> D&A <code>&lt;h4&gt;2020&lt;/h4&gt;</code> 2020	
b	굵게	<code>&lt;b&gt;D&amp;A&lt;/b&gt;</code> D&A	
u	밑줄	<code>&lt;u&gt;D&amp;A&lt;/u&gt;</code> D&A	
br	줄바꿈	<code>&lt;br&gt;D&amp;A&lt;/br&gt;</code> D&A	
a	링크 연결	<code>&lt;a href= 'http://naver.com' &gt;네이버&lt;/a&gt;</code> 네이버	href(연결할링크), title(확인문구)
span	문서 요소를 묶음	<code>&lt;span&gt;&lt;h5&gt;딥세션&lt;/h5&gt;&lt;/span&gt;</code> 딥세션	
li	목록의 리스트	<code>&lt;ul&gt;&lt;li&gt;2020&lt;/li&gt;&lt;/ul&gt;</code> 2020	
div	레이아웃	<code>&lt;div&gt;&lt;div&gt;&lt;/div&gt;&lt;/div&gt;</code>	
ul	순서가 없는 목록	<code>&lt;ul&gt;&lt;li&gt;&lt;/li&gt;&lt;/ul&gt;</code>	
ol	순서가 있는 목록	<code>&lt;ol&gt;&lt;li&gt;&lt;/li&gt;&lt;/ol&gt;</code>	
button	클릭할 버튼 생성	<code>&lt;button&gt;버튼&lt;/button&gt;</code>	type(종류)
p	한 단락	<code>&lt;p&gt;Deep Session 화이팅&lt;/p&gt;</code> Deep Session 파이팅	
img	이미지 파일 넣기	<code>&lt;img src= 'microphone.jpg' width= '100%' &gt;</code>	src(이미지주소), width(가로길이), height(세로길이)
input	입력도구 만들기	<code>&lt;input type= 'text' &gt;</code>	type(종류), value(초기값), placeholder(안내문)
textarea	큰 입력영역 만들기	<code>&lt;textarea rows= '5' cols= '10' &gt;&lt;/textarea&gt;</code>	rows(행수), cols(열수)

## 02 HTML의 구조

### 똑같은 태그가 엄청 많음

```
<!doctype html>
<html lang="ko" class="os_window">
  <head>...</head>
  <body style>
    <div id="daumIndex" class="d_index">...</div>
    <div id="daumWrap">
      <header id="daumHead" class="head_daum">...</header>
      <div id="tierContWrap" class="wrap_tiercont"></div>
      <div id="adLeft" class="wrap_tiercont" style="z-index:auto">...</div>
      <hr class="hide">
    <main id="daumContent">
      <div id="cMain" class="cont_main">
        <article id="mArticle" class="wrap_main">
          <div class="feature_tmp">
            <div id="adMain" class="advert_tmp">...</div>
            <div class="bg_login login_tmp #loginbox">...</div>
          </div>
          <div class="cmain_tmp">
            <div class="section_media">
              <h2 id="mediaTitle" class="screen_out">미디어</h2>
              <div id="mediaTab" data-tab="news" class="panel_bloc #newsbox news_on">...
            </div>
          </div> == $0
            <div id="kakaotv" class="section_multi" style="visibility: visible;" data-
            tiaramapper="kakaotv">...</div>
            <div class="section_blog" data-tiaramapper="story">...</div>
            <div class="section_channel" data-tiaramapper="channel">...</div>
          </div>
          <div class="wing_tmp">...</div>
          ::after
        </article>
      </div>
      <div id="cEtc" class="cont_topic" data-more=
      "news|377110,entertain|1003,entertain|50920,entertain|968,news|5014013" style=
      "visibility:hidden;">...</div>
    </main>
  </body>
</html>
```

※ 여러 개가 존재하는 태그에는  
대부분 class나 id로 구분함!!

id : #을 사용해 정보 찾기

class : . 을 사용해 정보 찾기

class는 여러 개가 존재할 수 있으며 각각의 클래스는 피어쓰기로 구분됨

※ 태그는 상하관계가 있기 때문에  
찾고자 하는 여러 태그가 존재한다면 이를 이용

자식 관계 : 바로 한단계 아래 하위태그, > 사용

자손관계 : 모든 하위태그, 피어쓰기 사용

선택자 > \*하면 모든 자식선택자 선택

# Requests

웹상의 데이터를 가져올 수 있는 패키지

```
!pip install requests  
import requests
```

# BeautifulSoup4

html코드를 파싱해 원하는 데이터를 추출할 수 있는 패키지

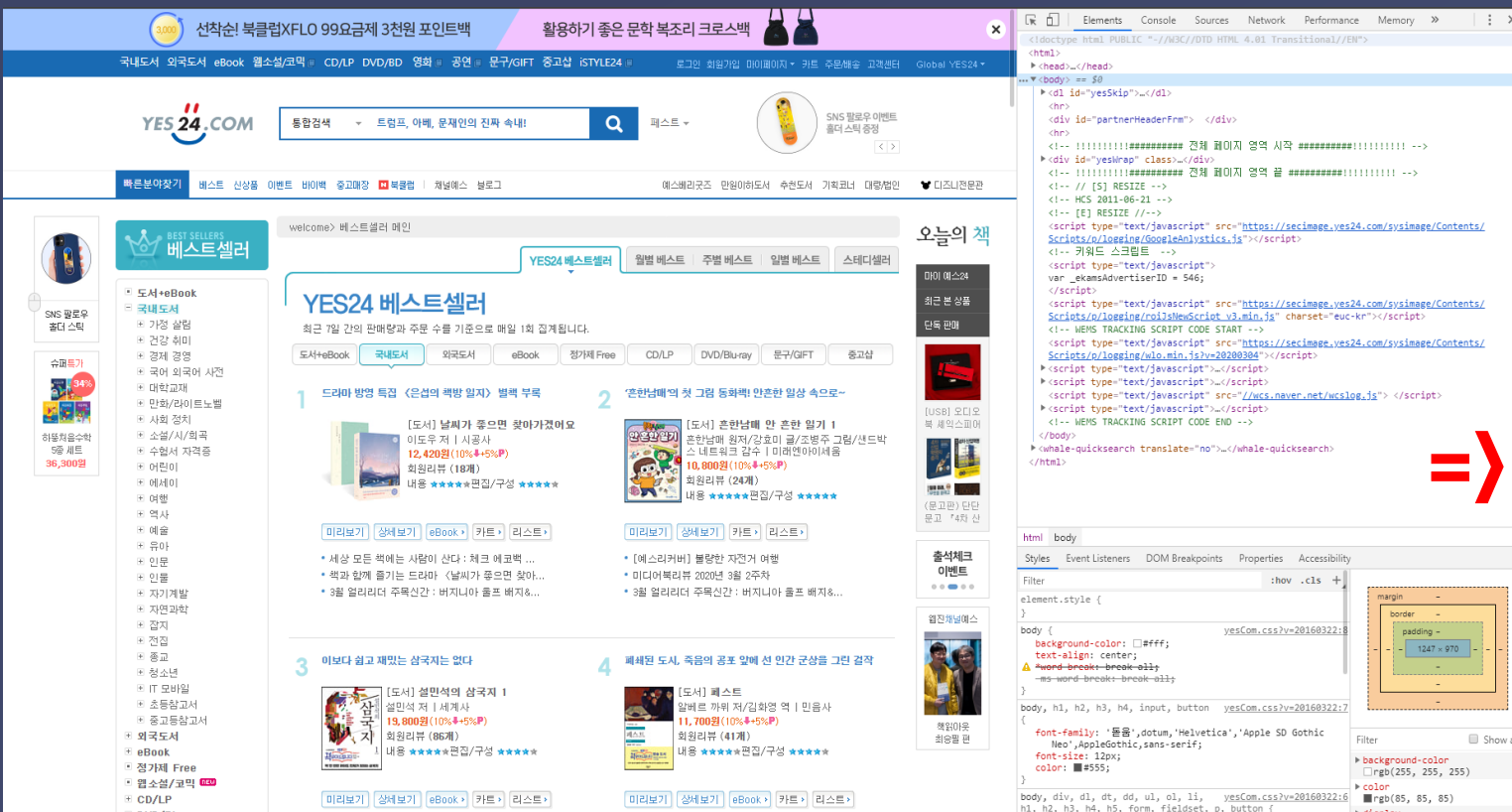
```
!pip install BeautifulSoup4  
from bs4 import BeautifulSoup
```

## 03 Requests와 BeautifulSoup4

```
In [2]: raw = requests.get('http://www.yes24.com/24/Category/BestSeller')
print(raw)
```

<Response [200]>

**get('웹페이지주소') : 웹페이지를 접속해 정보를 가져옴**



```
In [3]: raw.text
```

**text : 가져온 웹페이지의 소스코드**

```
Out [3]: 'rnrn <!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN">rnrn <html>rnrnrnrn<head><base href="http://www.yes24.com/24/" />rnrnrn<meta http-equiv="X-UA-Compatible" content="IE=5" />rnrnrnrn <meta http-equiv="Content-Type" content="text/html; charset=euc-kr" />rnrnrnrn<meta name="viewport" content="width=1170" />rnrnrnrn<title>YES24 | 대한민국 대표 인터넷서점 | 베스트셀러</title>rnrnrnrn<meta name="title" content="YES24 - 대한민국 대표 인터넷서점" />rnrn <meta name="description" content="YES24는 대한민국 1위 인터넷 온라인 서점입니다. 국내 최대의 도서정보를 보유하고 있으며, 음반, DVD, 공연, 영화까지 다양한 문화 콘텐츠 및 서비스를 제공합니다." />rnrn <meta name="keywords" content="인터넷 서점, 온라인 쇼핑, 상품 추천, 쇼핑물, 상품 검색, 도서 정보, 국내도서, 외국도서, 전자책, eBook, 이북, 크레마, 공연, 콘서트, 뮤지컬, 영화, 음반, 오디마, DVD, 블루레이, 예스24, YES24, 고보문고, 알라딘" />rnrn <meta property="og:image" content="https://secimage.yes24.com/sysimage/renew/logo_meta.png" />rnrnrnrn<script type="text/javascript" src="https://secimage.yes24.com/sysimage/Contents/Scripts/p/jquery/jquery-1.2.6.min.js"></script>rnrn<script type="text/javascript" src="https://secimage.yes24.com/sysimage/Contents/Scripts/p/jquery/jquery.menu-aim.js?v=20140801"></script>rnrn<script type="text/javascript" src="https://secimage.yes24.com/sysimage/Contents/Scripts/p/jquery/jquery.easing.1.3.min.js?v=20140801"></script>rnrn<script type="text/javascript" src="http://www.yes24.com/JavaScript/uti...
```

## 03 Requests와 BeautifulSoup4

```
In [4]: html = BeautifulSoup(raw.text, 'html.parser')
        print(html)
```

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN">
```

```
<html>
```

```
<head><base href="http://www.yes24.com/24/" />
```

```
<meta content="IE=5" http-equiv="X-UA-Compatible" />
```

```
<meta content="text/html; charset=utf-8" http-equiv="Content-Type" />
```

```
<meta content="width=1170" name="viewport" />
```

```
<title>YES24 | 대한민국 대표 인터넷서점 | 베스트셀러</title>
```

```
<meta content="YES24 - 대한민국 대표 인터넷서점" name="title" />
```

```
<meta content="YES24는 대한민국 1위 인터넷 온라인 서점 입니다. 국내 최대의 도서정보를 보유하고 있으며, 음반, DVD, 공연, 영화까지 다양한 문화 콘텐츠 및 서비스를 제공합니다." name="description" />
```

```
<meta content="인터넷 서점, 온라인 쇼핑, 상품 추천, 쇼핑물, 상품 검색, 도서 정보, 국내도서, 외국도서, 전자책, eBook, 이북, 크레마, 공연, 콘서트, 뮤지컬, 영화, 음반, 예매, DVD, 블루레이, 예스24, YES24, 교보문고, 알라딘" name="keywords" />
```

```
<meta content="https://secimage.yes24.com/sysimage/renew/logo_meta.png" property="og:image" />
```

```
<script src="https://secimage.yes24.com/sysimage/Contents/Scripts/p/jquery/jquery-1.2.6.min.js" type="text/javascript"></script>
```

```
<script src="https://secimage.yes24.com/sysimage/Contents/Scripts/p/jquery/jquery.menu-aim.js?v=20140801" type="text/javascript"></script>
```

```
<script src="https://secimage.yes24.com/sysimage/Contents/Scripts/p/jquery/jquery.easing.1.3.min.js?v=20140801" type="text/javascript">
```

**BeautifulSoup(웹페이지의 소스코드, 'html.parser')**  
**웹페이지의 소스코드를 html단위로 구분해줌**



### 03 Requests와 BeautifulSoup4

1. 컨테이너 구조를 파악 ※ 페이지 검사를 통해 필요한 정보가 담긴 컨테이너를 찾아야 함



2. 필요한 정보가 어떤 태그에 담겨져 있는지 확인



3. 다른 정보가 담겨져 있는 중복된 태그가 없는지 확인



4. 크롤러 만들기

드라마 방영 특집 <은섭의 책방 일지> 별책 부록

[도서] 날씨가 좋으면 찾아가겠어요  
이도우 저 | 시공사  
**12,420원 (10%↓+5%P)**  
회원리뷰 (18개)  
내용 ★★★★★편집/구성 ★★★★★

[미리보기](#) [상세보기](#) [eBook >](#) [카트](#) [리스트 >](#)

- 세상 모든 책에는 사람이 산다 : 체크 예코백 ...
- 책과 함께 즐기는 드라마 <날씨가 좋으면 찾아가...>
- 3월 얼리리더 주목시간 : 버지니아 울프 배지&...

1 드라마 방영 특집 <은섭의 책방 일지> 별책 부록

[도서] 날씨가 좋으면 찾아가겠어요  
이도우 저 | 시공사  
**12,420원 (10%↓+5%P)**  
회원리뷰 (18개)  
내용 ★★★★★편집/구성 ★★★★★

[미리보기](#) [상세보기](#) [eBook >](#) [카트](#) [리스트 >](#)

- 세상 모든 책에는 사람이 산다 : 체크 예코백 ...
- 책과 함께 즐기는 드라마 <날씨가 좋으면 찾아가...>
- 3월 얼리리더 주목시간 : 버지니아 울프 배지&...

2 '흔한남매'의 첫 그림 동화책! 안흔한 일상 속으로~

[도서] 흔한남매 안 흔한 일기 1  
흔한남매 원저/강효미 글/조병주 그림/샌드박스 네트워크 감수 | 미래엔아이세움  
**10,800원 (10%↓+5%P)**  
회원리뷰 (24개)  
내용 ★★★★★편집/구성 ★★★★★

[미리보기](#) [상세보기](#) [카트](#) [리스트 >](#)

- [예스리커버] 불량한 자전거 여행
- 미디어북리뷰 2020년 3월 2주차
- 3월 얼리리더 주목시간 : 버지니아 울프 배지&...

3 이보다 쉽고 재밌는 삼국지는 없다

[도서] 설민석의 삼국지 1  
설민석 저 | 세계사  
**19,800원 (10%↓+5%P)**  
회원리뷰 (86개)  
내용 ★★★★★편집/구성 ★★★★★

[미리보기](#) [상세보기](#) [eBook >](#) [카트](#) [리스트 >](#)

- YES24 X Disney 책 읽는 봄, 디즈니 스페셜! 나...
- 2019 올해의 책 선정 도서 24권을 공개합니다.
- <요즘 책방> 에서 소개한 책들

4 폐쇄된 도시, 죽음의 공포 앞에 선 인간 군상을 그린 걸작

[도서] 페스트  
알베르 카뮈 저/김화영 역 | 민음사  
**11,700원 (10%↓+5%P)**  
회원리뷰 (41개)  
내용 ★★★★★편집/구성 ★★★★★

[미리보기](#) [상세보기](#) [eBook >](#) [카트](#) [리스트 >](#)

- 민음사 세계문학 스프링 노트 증정!
- 세상 모든 책에는 사람이 산다 : 체크 예코백 ...
- YES24 X Disney 책 읽는 봄, 디즈니 스페셜! 나...

5 이보다 쉽고 재밌는 삼국지는 없다

[도서] 설민석의 삼국지 2  
설민석 저 | 세계사  
**19,800원 (10%↓+5%P)**  
회원리뷰 (43개)  
내용 ★★★★★편집/구성 ★★★★★

[미리보기](#) [상세보기](#) [eBook >](#) [카트](#) [리스트 >](#)

- <요즘 책방> 에서 소개한 책들

6 전승환이 꿈은 인생 문장

[도서] 내가 원하는 것을 나도 모를 때  
전승환 저 | 다산초당  
**14,400원 (10%↓+5%P)**  
회원리뷰 (94개)  
내용 ★★★★★편집/구성 ★★★★★

[미리보기](#) [상세보기](#) [eBook >](#) [카트](#) [리스트 >](#)

- 『내가 원하는 것을 나도 모를 때』 손글씨 공...
- 2020 신학기 인문교양 기획전
- 나를 찾는 인생의 문장을, 내가 원하는 것을 나...

### 03 Requests와 BeautifulSoup4



컨테이너의 태그  
`div#bestList ol > li`

홍보물 카피의 태그  
`div#bestList ol > p.copy`

가격의 태그  
`div#bestList ol > p.price`

```
In [5]: books = html.select('div#bestList ol > li')
print(books) select : html의 태그를 선택함 그 안의 모든 태그를 가져옴
```

```
[<li class="num1">
<p class="copy"><a href="/Product/Goods/86895523">드라마 방영 특집 <은섭의 책방
일지> 별책 부록</a></p>
<p class="image" id="location_0">
<a href="/Product/Goods/86895523">

</img></a>
</p>
<p>[도서] <a href="/Product/Goods/86895523">날씨가 좋으면 찾아가겠어요</a></p>
<p class="aupt"><a href="http://www.yes24.com//SearchCorner/Result?domain=ALL&
author_yn=Y&query=&auth_no=147966" target="_blank">이도우</a> 저 | <a href
="http://www.yes24.com//SearchCorner/Result?domain=ALL&company_yn=Y&query=
시공사">시공사</a></p>
<p class="price"><strong>12,420원</strong>(10%↓+5%P)
+5%
</p>
<p>회원리뷰 (<a href="/Product/Goods/86895523#Review">18개</a>)</p>
<p>
```

### 03 Requests와 BeautifulSoup4

```
In [10]: copy = books[0].select_one('p.copy')
print(copy)
```

**select\_one : html의 태그를 선택. 첫번째 하나의 태그만 가져옴**

<p class="copy"><a href="/Product/Goods/86895523">드라마 방영 특집 <은섭의 책방 일지> 별책 부록</a></p>

```
In [11]: copy = books[0].select_one('p.copy').text
print(copy)
```

드라마 방영 특집 <은섭의 책방 일지> 별책 부록

```
In [7]: price = books[0].select_one('p.price').text
print(price)
```

12,420원(10%+5%)

```
In [9]: for book in books:
        copy = book.select_one('p.copy').text
        price = book.select_one('p.price').text
        print(price, copy)
```

12,420원(10%+5%) 드라마 방영 특집 <은섭의 책방 일지> 별책 부록  
10,800원(10%+5%) '흔한남매'의 첫 그림 동화책! 안흔한 일상 속으로~  
19,800원(10%+5%) 이보다 쉽고 재미있는 삼국지는 없다  
11,700원(10%+5%) 폐쇄된 도시, 죽음의 공포 앞에 선 인간 군상을 그린 걸작  
19,800원(10%+5%) 이보다 쉽고 재미있는 삼국지는 없다  
14,400원(10%+5%) 전승환이 꿈은 인생 문장  
14,400원(10%+5%) 하버드 성공 공식 "신호를 차단하고 깊이 몰입하라"  
14,400원(10%+5%) 하루 1분이면 교양 시민  
24,120원(10%+5%) 리더는 타고나는가 만들어지는가  
10,800원(10%+5%) 설렘과 함께 위인들의 어린 시절을 만나요!  
12,150원(10%+5%) 현실에서 딱 1cm 벗어나 행복찾기!  
14,850원(10%+5%) 영화 <작은 아씨들> 공식 오리지널 커버 디자인  
10,800원(10%+5%) 흔한 일상 속에 숨겨진 반짝거리는 웃음과 재미!  
11,610원(10%+5%) 2020 최신판으로 토익 기출 보카 업그레이드  
17,820원(10%+5%) 지대넓얕 0권, 모든 지식의 시작  
16,650원(10%+5%) 우리는 코스모스에서 태어났다! 인류를 품은 거대한 우주, 그 기원과 신비, 인간의 본질을 밝히다  
11,700원(10%+5%) '다시 한 번 사랑해보기로 한' 모든 이를 향한 따뜻한 시선  
10,800원(10%+5%) 고통과 공간의 능력을 깨우치게 하는 소설  
25,200원(10%+5%) 서울대 도서관 대출도서 부동의 1위! 풀리처상 수상작품  
14,400원(10%+5%) 푸른 길잡이, 자연 경관지도, 여행하는 보석, 여행

이렇게 for문을 통해  
여러 개의 데이터를  
한 번에 수집할 수 있음!

## 03 Requests와 BeautifulSoup4

### 불러온 웹페이지에서 다른 페이지로 넘어가 정보를 가져올 때

```
In [30]: raw2 = requests.get('http://yes24.com'+books[0].select_one('p.copy a').attrs['href'])

In [32]: html2 = BeautifulSoup(raw2.text, 'html.parser')
         title = html2.select_one('h2.gd_name').text
         print(title)

날씨가 좋으면 찾아가겠어요
```

특정 태그에 속해있는 href태그를 불러오면  
연결된 페이지의 정보가 불러와짐

### 이미지 정보를 가져올 때

```
In [33]: from urllib.request import urlopen

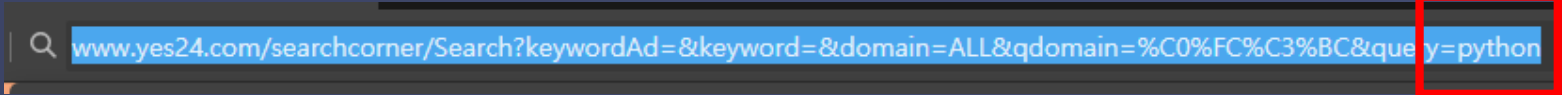
In [36]: img = books[0].select_one('ol li p.image img').attrs['src']

In [39]: urlopen(img, '책1.png')
```

src로 불러와짐

# 03 Requests와 BeautifulSoup4

## ※ 검색 자동화



YES24.COM 통합검색 python 김새해 추천책

특은분야찾기 베스트 신상품 이벤트 바이백 중고매장 북클럽 채널메스 블로그 에스베리굿즈 만원이하도

통합검색 4,026

국내도서 490  
외국도서 2,606  
eBook 301  
CD/LP 16  
DVD/BD 42  
영화 4  
문구/GIFT 6  
중고샵 428  
리뷰 127  
기사, 인터뷰 6

결과내 재검색  
입력후 엔터

책소개 검색  
입력후 엔터

목차 검색 (더보기)  
2020 시나공 정보처...  
2020 시나공 정보처...  
파이썬으로 배우는 ...  
실습으로 배우는 S...  
Microsoft Power BI...

상세 검색 >

국내도서 | 외국도서 | eBook | CD/LP | DVD/BD | 문구/GIFT | 중고샵 | 영화

- IT 모바일 (469)
- 어린이 (20)
- 대학교재 (19)
- 자연과학 (5)
- 경제 경영 (5)
- 사회 정치 (1)
- 전집 (1)

[작가] Python Lee Jackson (~) 피톤 리 잭슨  
최신작 [CD] Python Lee Jackson - In A Broken Dream (LP Miniature)  
동명미인 작가보기 +

통합검색 : "python" 검색결과 1-20 / 3,461건

인기도 | 정확도 | 신상품 | 최저가 | 최고가 | 평점순 | 리뷰순 20개 옵션선택 품질포함

1 [도서] 혼자 공부하는 파이썬 : 파이썬 최신 버전-혼자 공부하는 시리즈  
윤인성 저 | 한빛미디어 | 2019년 06월  
16,200원 (10% 할인) 900원  
판매지수 72,960 | 회원리뷰(31건) ★★★★★ 9.4  
출고 예상일 : 1 일 이내 안내  
#2019IT연말결산선정도서 #새해결심 #분월 #کم공 #프로그래머  
사은품 『혼자 공부하는 시리즈』 혼공노트 증정 이벤트  
사은품 기획권 [대학생/취준생] 신학기/스펙/취업 혜택 총 집합!  
기획권 파이썬, 나도 한번 해볼까?  
본월서비스 이용이 가능한 도서입니다. 자세히보기 >

2 [도서] Do it! 점프 투 파이썬 (파이썬 3 최신 버전 반영, 개정판)  
박용홍 저 | 이지스퍼블리싱 | 2019년 06월  
16,920원 (10% 할인) 15,228원

```
In [16]: keyword = 'python'

In [17]: raw = requests.get('http://www.yes24.com/searchcorner/Search?keywordAd=&keyword=&domain=ALL&qdomain=%C0%FC%C3%BC&query='+keyword)
```

## for문을 통해 자동화 가능

```
In [19]: keyword = ['python', 'java', 'r']

In [28]: for k in keyword:
raw = requests.get('http://www.yes24.com/searchcorner/Search?keywordAd=&keyword=&domain=ALL&qdomain=%C0%FC%C3%BC&query='+k)
html = BeautifulSoup(raw.text, 'html.parser')
books = html.select('td.goods_infgrp')
for book in books[:20]:
title = book.select_one('a').text
writer = book.select_one('div.goods_info a').text
price = book.select_one('div.goods_price em.yes_b').text
print(title, writer, price)
```

혼자 공부하는 파이썬 윤인성 16,200  
Do it! 점프 투 파이썬 박용홍 16,920  
밑바닥부터 시작하는 딥러닝 사이토 고키 21,600  
파이썬 머신러닝 완벽 가이드 권철민 34,200  
밑바닥부터 시작하는 딥러닝 2 사이토 고키 26,100  
파이썬과 리엑트를 활용한 주식 자동거래 시스템 구축 박재현 25,200  
케라스 창시자에게 배우는 딥러닝 프랑소와 술레 29,700  
파이썬으로 배우는 프로그래밍 기초 크리스 로피 13,000  
머신 러닝 교과서 with 파이썬, 사이킷런, 텐서플로 세바스찬 라시카 29,700  
파이썬 기초 문법 이원하 0  
파이썬 Python 김영천 25,000  
파이썬 연습 김기용 26,000  
알기 쉬운 파이썬 SQL 코딩하기 정현희 13,500  
파이썬으로 배우는 수치 데이터 처리 김홍근 19,800  
딥러닝 AI 프로젝트 실사례 Santanu Pattanayak 29,000  
예제로 배우는 자연어 처리 기초 소홀 고시 31,500  
모두의 파이썬 이승찬 10,800  
Do it! 점프 투 파이썬 + 딥러닝 입문 박용홍 34,740  
실무자를 위한 파이썬 Python 100제 오승환 18,000  
모두의 데이터 분석 with 파이썬 송석리 16,200  
글자 바로 쓰기 초등 국어 1-1 (2020년용) 미래엔 초등 국어 연구회 7,650  
하루 한장 국어 + 수학 1-1 세트 (2020년용) 미래엔콘텐츠연구회 26,460  
글자 바로 쓰기 초등 국어 2-1 (2020년용) 미래엔 초등 국어 연구회 7,650



### 03 Requests와 BeautifulSoup4

#### ※ 페이지 넘기기 자동화

www.yes24.com/searchcorner/Search?keywordAd=&keyword=&domain=ALL&qdomain=%C0%FC%C3%BC&query=python

www.yes24.com/searchcorner/Search?keywordAd=&keyword=&domain=ALL&qdomain=%c0%fc%c3%bc&query=python&PageNumber=2&scode=012

## 아까와 마찬가지로 페이지 주소를 살펴보면 답이 나옴!

```
In [32]: for k in keyword:
          for num in range(1,10):
            raw = requests.get(
              'http://www.yes24.com/searchcorner/Search?keywordAd=&keyword=&domain=ALL&qdomain=%C0%FC%C3%BC&query='+k
             +'&PageNumber='+str(num))
            html = BeautifulSoup(raw.text, 'html.parser')
            books = html.select('td.goods_infogrp')
            for book in books[:20]:
              title = book.select_one('a').text
              writer = book.select_one('div.goods_info a').text
              price = book.select_one('div.goods_price em.yes_b').text
              print(title,writer,price)
```

혼자 공부하는 파이썬 윤인성 16,200  
Do it! 점프 투 파이썬 박응용 16,920  
말바닥부터 시작하는 딥러닝 사이토 고키 21,600  
파이썬 머신러닝 완벽 가이드 권철민 34,200  
말바닥부터 시작하는 딥러닝 2 사이토 고키 26,100  
파이썬과 리액트를 활용한 주식 자동거래 시스템 구축 박재현 25,200  
케라스 창시자에게 배우는 딥러닝 프랑소와 솔레 29,700  
파이썬으로 배우는 프로그래밍 기초 크리스 로피 13,000  
머신 러닝 교과서 with 파이썬, 사이킷런, 텐서플로 세바스찬 라스카 29,700  
파이썬 기초 문법 이원하 0  
파이썬 Python 김영천 25,000  
파이썬 연습 김기용 26,000  
알기 쉬운 파이썬 SQL 코딩하기 정현희 13,500  
파이썬으로 배우는 수치 데이터 처리 김동근 19,800  
딥러닝 AI프로젝트 실사례 Santanu Pattanayak 29,000  
예제로 배우는 자연어 처리 기초 쇼훔 고시 31,500  
모두의 파이썬 이승찬 10,800  
Do it! 점프 투 파이썬 + 딥러닝 입문 박응용 34,740  
실무자를 위한 파이썬 Python 100제 오승환 18,000  
모두의 파이썬 분석 파이썬 파이썬 추천가 16,200

### 03 Requests와 BeautifulSoup4

※ 크롤링하려고 하던 페이지가 막혀서 접근불가일 때  
해결방안 : headers={'User-Agent': '접근자'} 를 설정

```
In [33]: raw = requests.get('https://movie.naver.com/movie/running/current.nhn', headers={'User-Agent': 'Mozilla/5.0'})  
html = BeautifulSoup(raw.text, 'html.parser')
```

Mozilla/5.0의 탈을 쓰고 웹페이지에 접근한다는 뜻

※ 똑같은 이름의 태그라 원하는 정보를 수집하지 못할 때

```
In [36]: movies = html.select('dl.lst_dsc')  
  
for movie in movies:  
    title = movie.select_one('dt.tit_a').text  
    genre = movie.select_one('dl.info_txt1 dd a').text  
    actor = movie.select_one('dl.info_txt1 dd a').text  
    print(title, genre, actor)
```

인비저블맨 공포 공포  
다크 워터스 드라마 드라마  
1917 드라마 드라마  
더 보이 2: 돌아온 브람스 공포 공포  
지푸라기라도 잡고 싶은 짐승들 범죄 범죄  
작은 아씨들 드라마 드라마  
정직한 후보 코미디 코미디  
시원찮은 그녀를 위한 육성방법 피날레 애니메이션 애니메이션  
찬실이는 복도 많지 드라마 드라마  
스타 이즈 본 드라마 드라마  
울프 콜 액션 액션  
악몽 미스터리 미스터리  
첼렌저 범죄 범죄  
타오르는 여인의 초상 드라마 드라마  
리암 갤러거 다큐멘터리 다큐멘터리  
작가 미상 드라마 드라마  
하이, 찍시 코미디 코미디  
슬럼독 밀리어네어 범죄 범죄  
비긴 어게인 드라마 드라마  
유리알에 갇힌 자들 미스터리 미스터리

```
<dl class="info_txt1">  
  <dt class="tit_t1">개요</dt>  
  <dd> == $0  
    <span class="link_txt">  
      <a href="/movie/sdb/browsing/bmovie.nhn?genre=1">드라마</a>  
    <!-- N=a.no1.genre,r:1 -->  
  </span>  
</dd>  
</dl>
```

```
<dl class="info_txt1">  
  <dt class="tit_t1">개요</dt>  
  <dd>...</dd>  
  <dt class="tit_t2">감독</dt>  
  <dd>...</dd>  
  <dt class="tit_t3">출연</dt>  
  <dd>  
    <span class="link_txt"> == $0  
      <a href="/movie/bi/pi/basic.nhn?code=6176">마크 러팔로</a>  
    <!-- N=a.no1.actor,r:1 -->  
  </span>  
</dd>  
</dl>
```

## 03 Requests와 BeautifulSoup4

해결방안 : 같은 태그 뒤에 :nth-of-type(태그의 순서)를 덧붙임

```
In [37]: movies = html.select('dl.lst_dsc')
```

```
for movie in movies:
```

```
    title = movie.select_one('dt.tit a').text
```

```
    genre = movie.select_one('dl.info_txt1 dd:nth-of-type(1) a').text
```

```
    actor = movie.select_one('dl.info_txt1 dd:nth-of-type(3) a').text
```

```
    print(title, genre, actor)
```

info\_txt1 다음의 dd가 여러 개 있는 게 겹치는 것이기 때문에  
dd:nth-of-type사용

인비저블맨 공포 엘리자베스 모스

다크 워터스 드라마 마크 러팔로

1917 드라마 조지 맥케이

더 보이 2: 돌아온 브람스 공포 케이티 홈즈

지푸라기라도 잡고 싶은 짐승들 범죄 전도연

작은 아씨들 드라마 시얼샤 로넌

정직한 후보 코미디 라미란

시원찮은 그녀를 위한 육성방법 피날레 애니메이션 마츠오카 요시츠구

찬슬이는 복도 많지 드라마 강말금

스타 이즈 본 드라마 브래들리 쿠퍼

울프 콜 액션 프랑수아 시빌

악몽 미스터리 오지호

젠틀맨 범죄 매튜 맥커너히

타오르는 여인의 초상 드라마 아델 하에넬

리암 갤러거 다큐멘터리 리암 갤러거

작가 미상 드라마 톰 윌링

하이, 잭시 코미디 아담 드바인

슬럼독 밀리어네어 범죄 데브 파텔

비긴 어게인 드라마 키이라 나이틀리

유기체 공화국 리즈 리스



### ※ 정적페이지와 동적페이지?

정적페이지는 웹페이지의 화면이  
클릭이나 입력을 추가했을 때 웹페이지 주소가 변경되는 것

동적페이지는 웹페이지 주소가 변경되지는 않지만  
보이는 웹페이지가 변경되는 페이지

javascript는 외관적인 페이지를 변경해주는 용도이기 때문에 동적페이지 수집방법으로 수집

## 04 Selenium

# Selenium

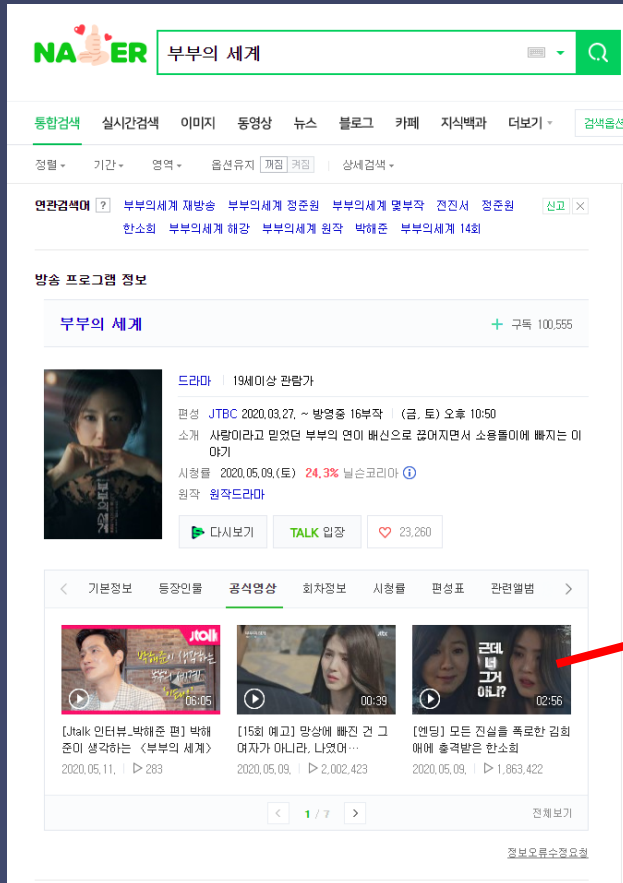
웹 브라우저를 컨트롤하여 웹을 자동화시켜주는 패키지

`!pip install selenium`

웹자동화를 시키기 위한 크롬드라이버를 먼저 설치해야함  
크롬 실행후 도움말에서 Chrome정보 클릭후 **버전확인** 후 **다운받기**  
(<http://chromedriver.chromium.org/downloads>)

## 04 Selenium

# 클릭을해서 정보를 추가로 수집할 수 있음



NAVER 부부의 세계

통합검색 실시간검색 이미지 동영상 뉴스 블로그 카페 지식백과 더보기

정렬 기간 영역 옵션유지 캐집 캐집 상세검색

연관검색어 부부의세계 재방송 부부의세계 정준원 부부의세계 딸부자 전진서 정준원 한소희 부부의세계 해강 부부의세계 원작 박해준 부부의세계 14회

방송 프로그램 정보

부부의 세계 + 구독 100,555

드라마 | 19세이상 관람가

편성 JTBC 2020.03.27. ~ 방영종 16부작 | (금, 토) 오후 10:50

소개 사랑이라고 믿었던 부부의 연이 배신으로 끊어지면서 소용돌이에 빠지는 이야기

시청률 2020.05.09.(토) 24.3% 날순코리아

원작 원작드라마

다시보기 TALK 입장 23,260

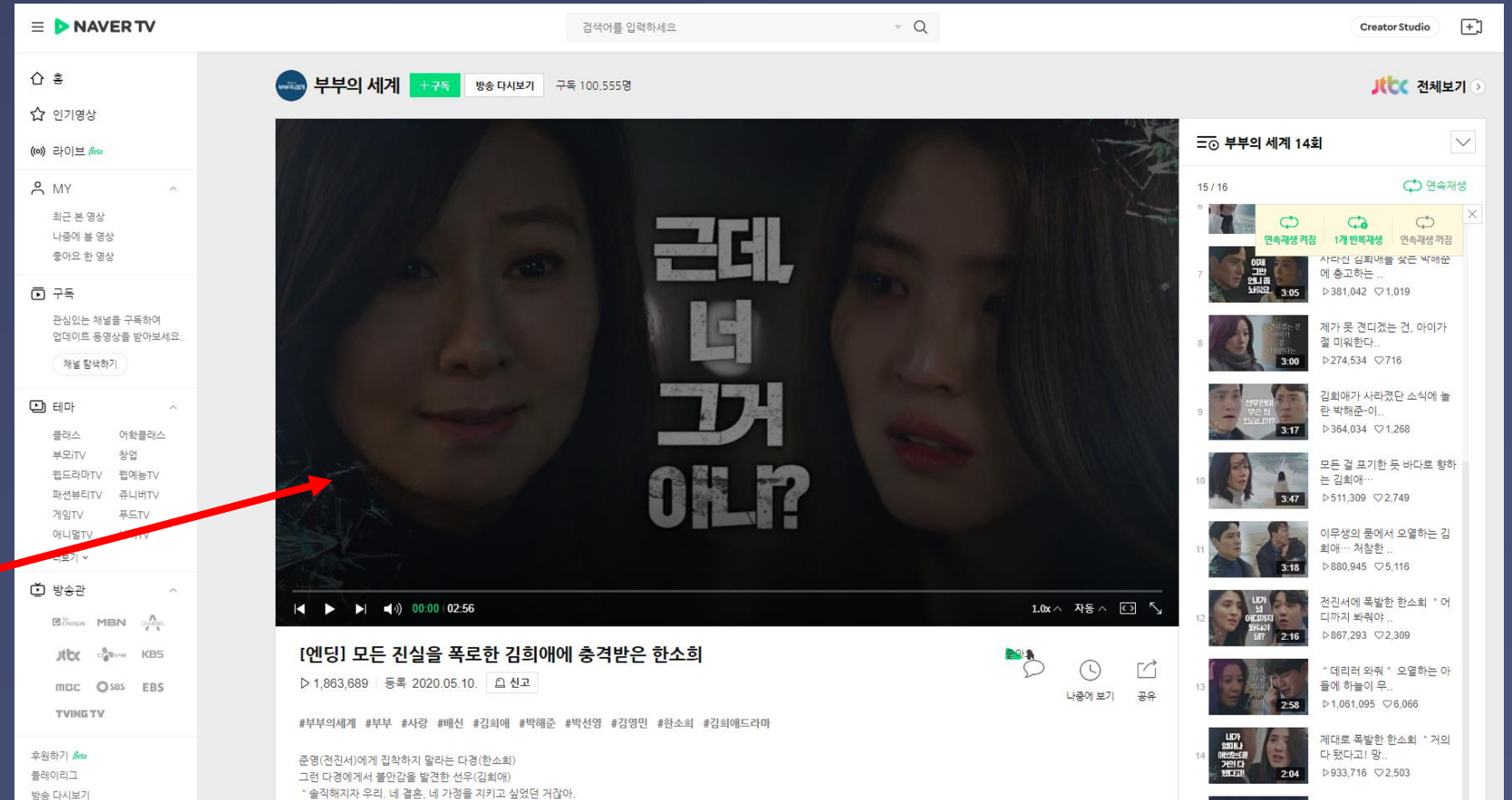
기본정보 등장인물 공식영상 회차정보 시청률 편성표 관련앨범

[15회 예고] 망상에 빠진 건 그 여자가 아니라, 나였어...

[연딩] 모든 진실을 폭로한 김희애에 충격받은 한소희

전체보기

정보오픈소셜유포



NAVER TV 검색어를 입력하세요

부부의 세계 +구독 방송 다시보기 구독 100,555명

전체보기

부부의 세계 14회

15 / 16

연속재생

연속재생 커짐 1개 반복재생 연속재생 커짐

사다인 김희애를 앓은 박해준에 충고하는...

381,042 3.05

제가 못 견디겠는 건, 아이가 결 미워한다.

274,534 3.00

김희애가 사라졌단 소식이 놀란 박해준-이..

364,034 3.17

모든 걸 포기한 듯 바다로 향하는 김희애...

511,309 3.47

이무생의 품에서 오열하는 김희애... 처절한...

880,945 3.18

전진서에 폭발한 한소희 "어디까지 봐줘야..."

867,293 2.16

"데리러 와줘" 오열하는 아들이 하늘이 무..

1,061,095 2.58

게대로 폭발한 한소희 "거의 다 왔다고! 망.."

933,716 2.04

근데 너 그거 아니?

[엔딩] 모든 진실을 폭로한 김희애에 충격받은 한소희

1,863,689 등록 2020.05.10. 신고

#부부의세계 #부부 #사랑 #배신 #김희애 #박해준 #박선영 #김영민 #한소희 #김희애드라마

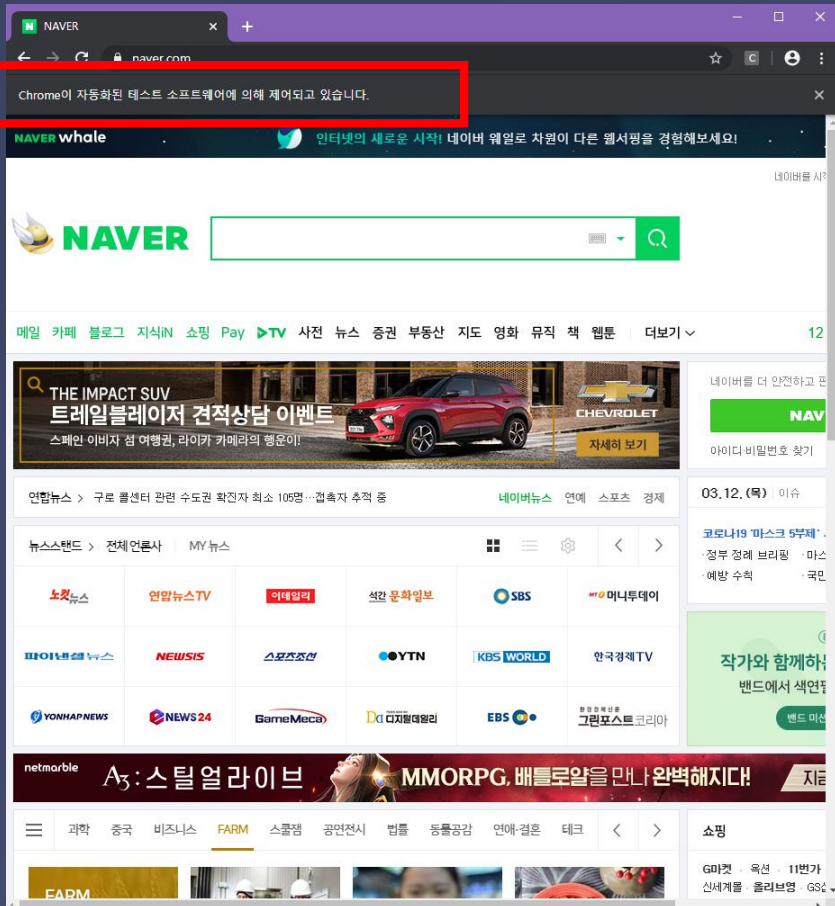
준영(전진서)에게 집착하지 말라는 다경(한소희) 그런 다경에게서 불안감을 발견한 선우(김희애) "술 취하자 우리. 네 결혼, 네 가정을 지키고 싶었던 거잖아."

후원하기 Auto 클라이그 방송 다시보기

## 04 Selenium

```
In [48]: from selenium import webdriver
```

```
In [52]: driver = webdriver.Chrome('./chromedriver')  
driver.get('https://www.naver.com')  
driver.close()
```



실행만 시키면 크롬창이 실행되며  
원하는 페이지로 자동으로 이동

※ 웹드라이버는 페이지를 로드, 이동하는데  
시간이 걸리기 때문에 성급하게 페이지를 닫기보다는  
기다리세요!  
(로드까지 1~3초정도의 시간이 소요됨  
네트워크환경에 따라 다를 수 있음)

## 04 Selenium

**send\_keys() : input 태그에 값을 입력해주는 메소드**

```
In [50]: keyword = driver.find_element_by_css_selector('span.green_window input')  
keyword.send_keys('부부의 세계')
```

**click() : button 태그를 클릭해주는 메소드**

```
In [67]: search = driver.find_element_by_css_selector('button#search_btn')  
search.click()
```



## 04 Selenium

### Requests와 마찬가지로 for문을 사용해 자동화를 시킬 수 있음

```
In [53]: video = driver.find_elements_by_css_selector('div.formula_video li')
for v in video:
    title = v.find_element_by_css_selector('strong a').text
    date = v.find_element_by_css_selector('div.bt_desc span.update').text
    view = v.find_element_by_css_selector('div.bt_desc span.play_count').text
    print(title, date, view)
```

[Jtalk 인터뷰\_박해준 편] 박해준이 생각하는 〈부부의 세계〉 '이태오' 2020.05.11. 재생수 9,150  
[15회 예고] 망상에 빠진 건 그 여자가 아니라, 나였어... 2020.05.09. 재생수 2,070,949  
[엔딩] 모든 진실을 폭로한 김희애에 충격받은 한소희 2020.05.09. 재생수 1,908,645

※ 간혹 나는 잘 찾았는데 선택자를 찾을 수 없다는 에러가 발생

대표원인 : 대부분 드라이버가 실행되는 시간 보다  
파이썬의 코드가 읽히는 시간이 더 짧기 때문

해결방안 :

```
import time
time.sleep(5)
```

## 04 Selenium

### 데이터가 없을 때 오류가 발생할 수 있음

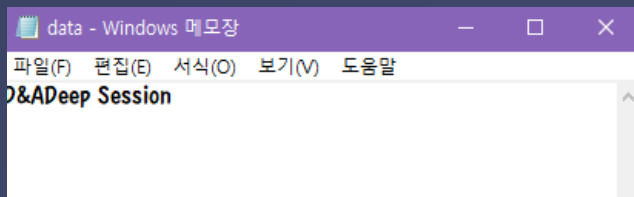
```
In [79]: video = driver.find_elements_by_css_selector('div.formula_video li')
         for v in video:
             title = v.find_element_by_css_selector('strong a').text
             date = v.find_element_by_css_selector('div.bt_desc span.update').text
             try:
                 view = v.find_element_by_css_selector('div.bt_desc span.play_count').text
             except:
                 view = '정보없음'
             print(title,date,view)
```

=> 예외처리(`try, except`)를 사용하여 에러 해결

## 05 데이터 저장하기

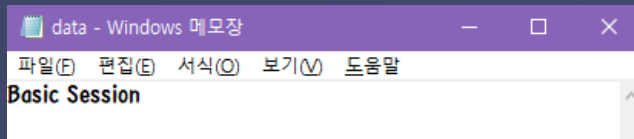
'w' : 쓰기모드  
'a' : 추가모드  
'r' : 읽기모드

```
In [80]: f = open('data.txt', 'w')  
f.write('D&A')  
f.write('Deep Session')  
f.close()
```



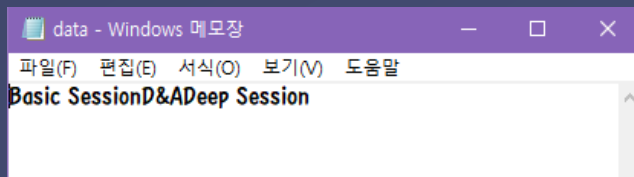
write를 사용하면 원래 데이터에  
이어서 작성됨

```
In [89]: f = open('data.txt', 'w')  
f.write('Basic Session')  
f.close()
```



다시 데이터를 오픈해서 write를 하면  
이전 데이터는 사라지고 덮어쓰기 함

```
In [90]: f = open('data.txt', 'a')  
f.write('D&A')  
f.write('Deep Session')  
f.close()
```



a 모드를 사용해야  
내용을 추가할 수 있음

```
In [91]: f = open('data.txt', 'r')  
data = f.read()  
print(data)  
f.close()
```

Basic SessionD&A\nDeep Session

읽기모드는 쓰거나 추가할 수는 없음



## 05 데이터 저장하기

```
In [103]: f = open('data.csv', 'w')  
f.write('학회원1,학회원2,학회원3')  
f.write('\n공영경,김보현,신지섭')  
f.write('\n권오현,오주영,최홍록')  
f.close()
```

```
In [104]: f = open('data.csv', 'r')  
data = f.read()  
print(data)  
f.close()
```

학회원1,학회원2,학회원3  
공영경,김보현,신지섭  
권오현,오주영,최홍록

자동 저장

파일 홈 삽입 페이지 레이아웃 수식 데이터 검토 보기 도움말 검색

잘라내기 붙여넣기 복사 서식 복사 클립보드 글꼴 맞춤

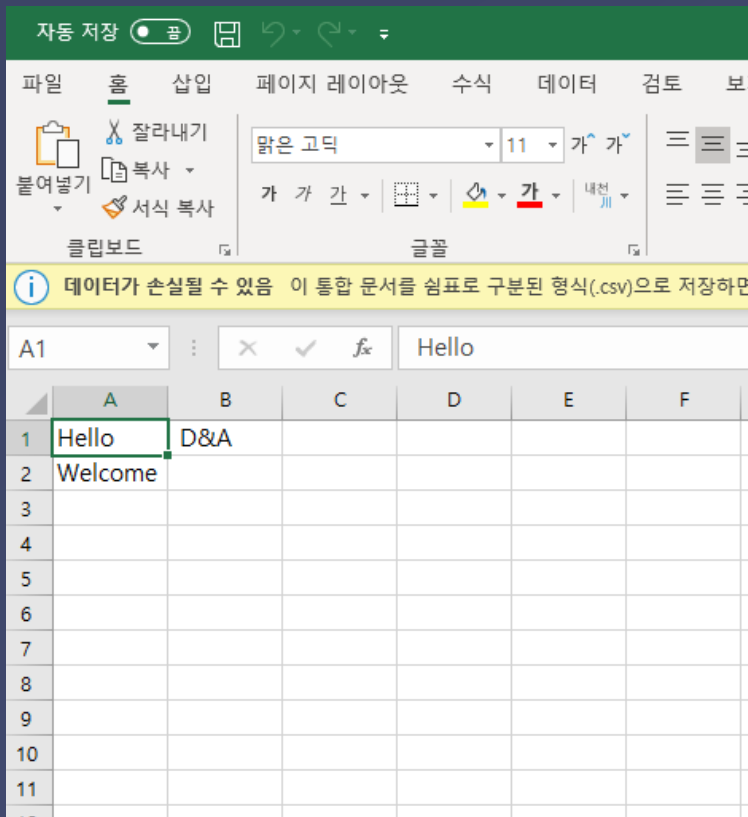
데이터가 손실될 수 있음 이 통합 문서를 심표로 구분된 형식(.csv)으로 저장하면 일부 기능이 손실될 수 있습니다

	A	B	C	D	E	F	G	H	I
1	학회원1	학회원2	학회원3						
2	공영경	김보현	신지섭						
3	권오현	오주영	최홍록						
4									
5									
6									
7									
8									
9									
10									
11									
12									
13									
14									
15									

## 05 데이터 저장하기

```
In [108]: a1 = 'Hello, D&A'
a2 = '\nWelcome'
f = open('data.csv', 'w')
f.write(a1)
f.write(a2)
f.close()
```

csv로 저장하면 콤마때문에  
a1 셀이 나뉘지게 된다.



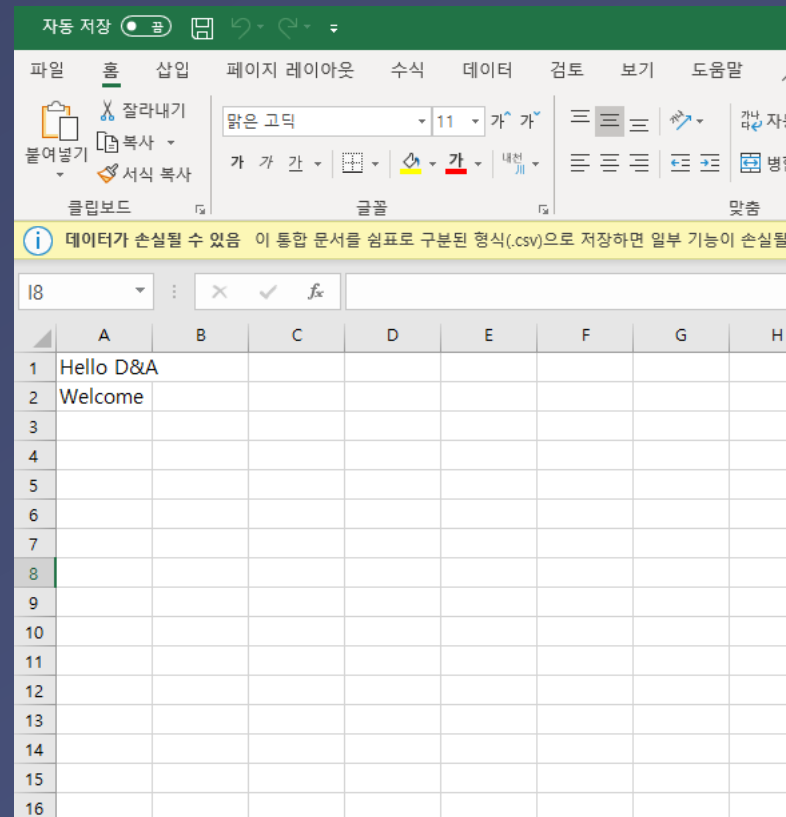
자동 저장 (On) | 파일 | 홈 | 삽입 | 페이지 레이아웃 | 수식 | 데이터 | 검토 | 보기

데이터가 손실될 수 있음 이 통합 문서를 심표로 구분된 형식(.csv)으로 저장하면

	A	B	C	D	E	F
1	Hello	D&A				
2	Welcome					
3						
4						
5						
6						
7						
8						
9						
10						
11						

```
In [110]: a1 = 'Hello, D&A'
a2 = '\nWelcome'
f = open('data.csv', 'w')
f.write(a1.replace(',', ''))
f.write(a2)
f.close()
```

replace를 통해 교체를  
해줘야함



자동 저장 (On) | 파일 | 홈 | 삽입 | 페이지 레이아웃 | 수식 | 데이터 | 검토 | 보기 | 도움말

데이터가 손실될 수 있음 이 통합 문서를 심표로 구분된 형식(.csv)으로 저장하면 일부 기능이 손실될

	A	B	C	D	E	F	G	H
1	Hello D&A							
2	Welcome							
3								
4								
5								
6								
7								
8								
9								
10								
11								
12								
13								
14								
15								
16								

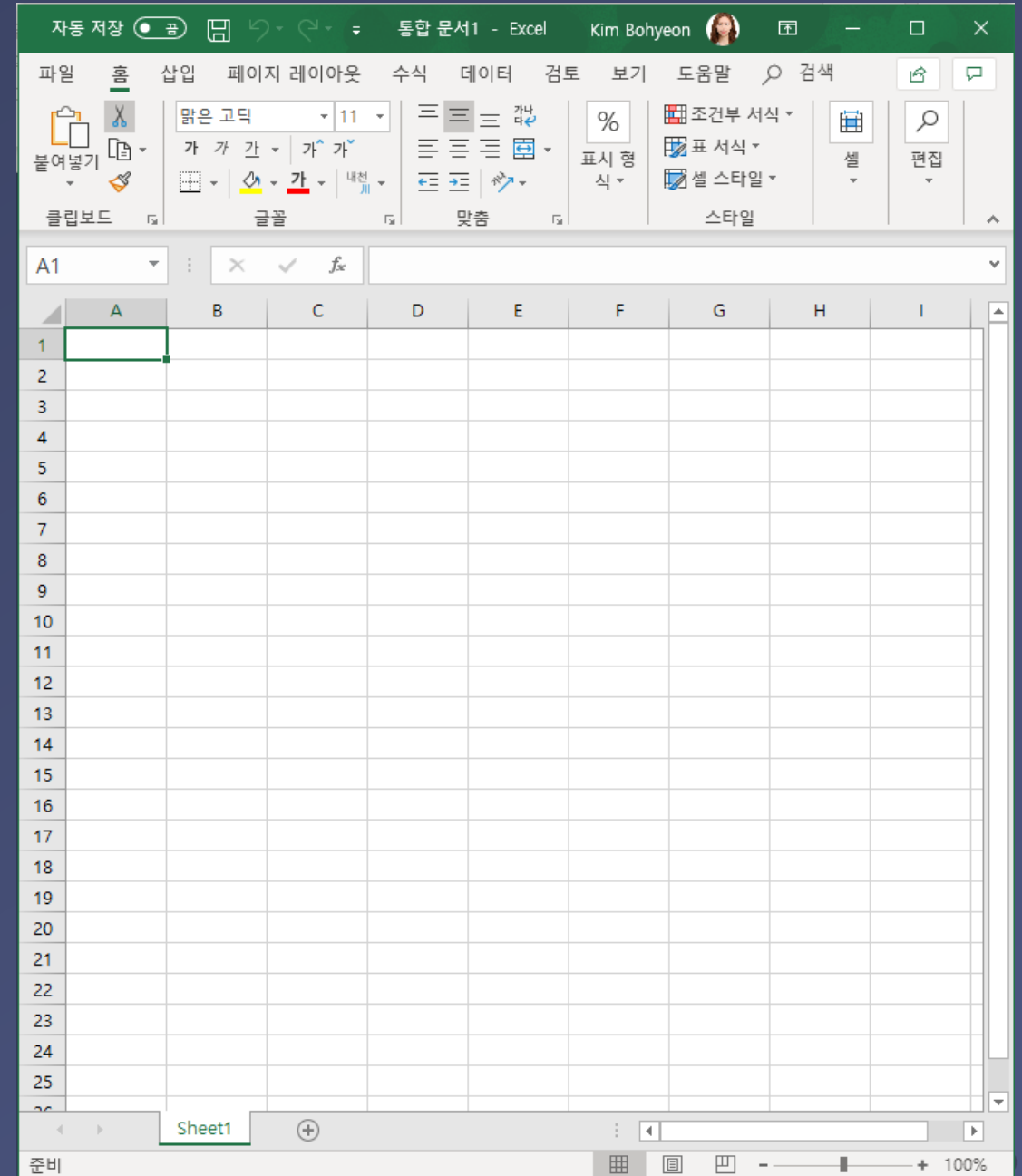
## 05 데이터 저장하기

```
In [112]: import openpyxl

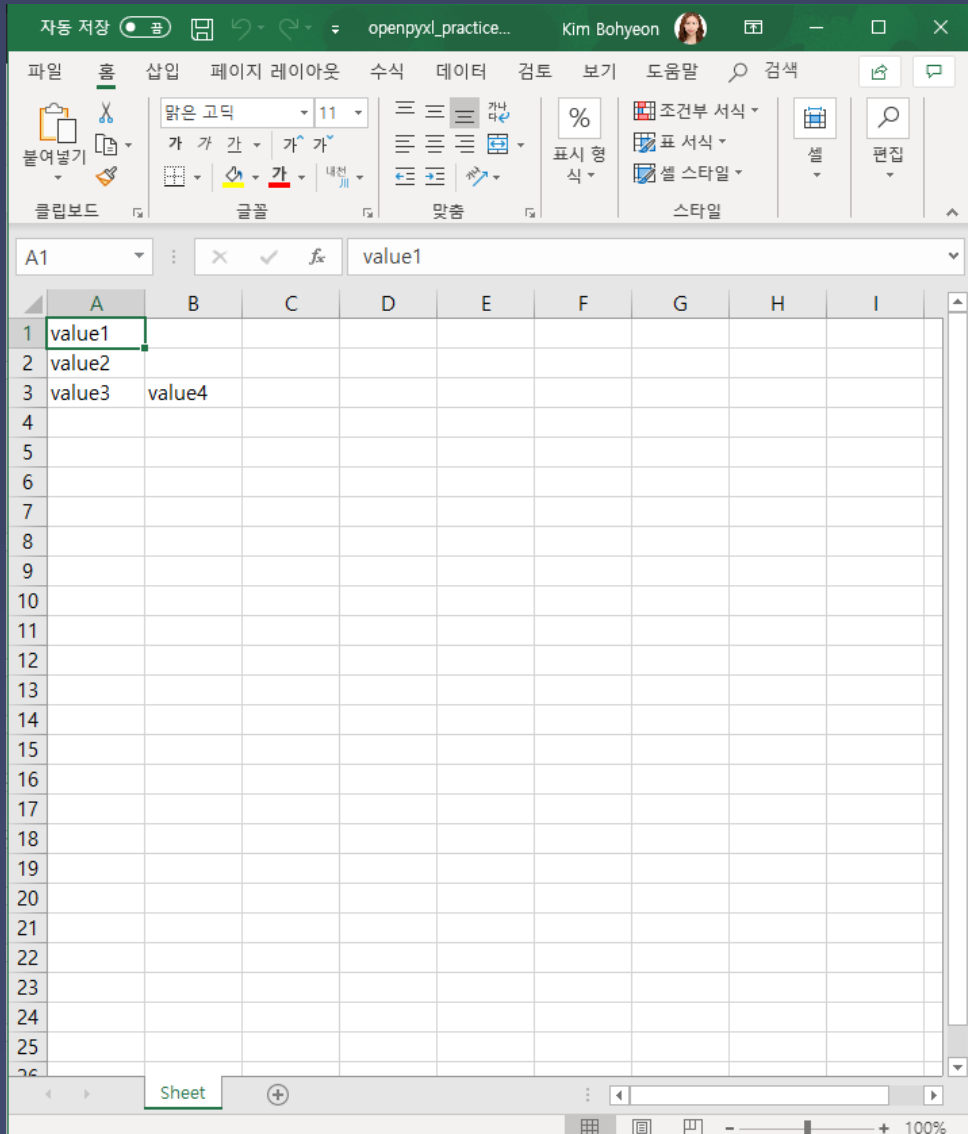
wb = openpyxl.Workbook()
sheet = wb.active
```

**openpyxl.Workbook()**  
새로운 엑셀 워크북을 열기

**wb.active**  
엑셀 워크시트를 활성화시켜  
데이터를 넣을 수 있게 함



## 05 데이터 저장하기



```
In [132]: sheet['A1'] = 'value1'
```

```
In [133]: sheet.cell(row=2,column=1).value = 'value2'
```

```
In [134]: sheet.append(['value3','value4'])
```

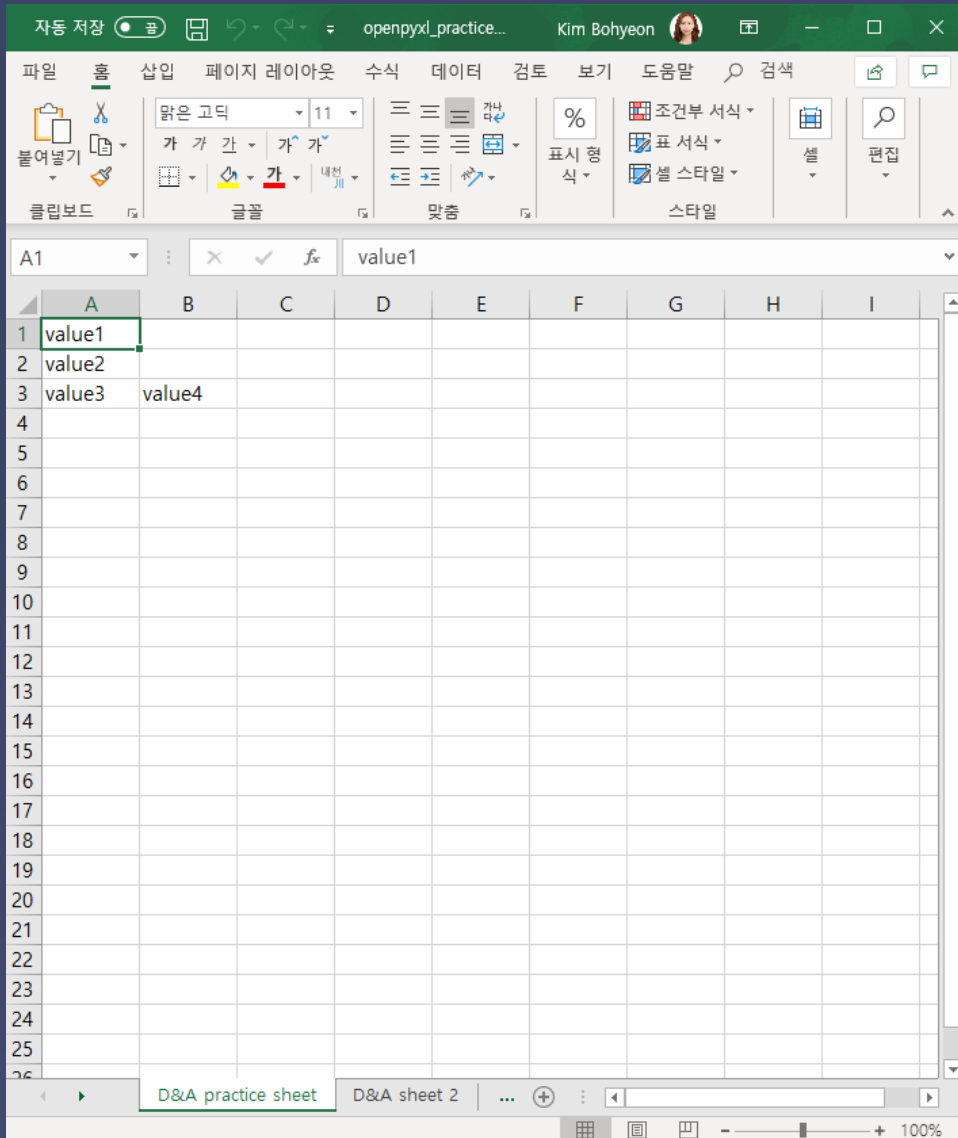
## 값을 채우는 여러 가지 방법

`sheet['셀이름'] = 값`

`sheet.cell(row=n,column=n).value = 값`

`sheet.append(['값'])`

## 05 데이터 저장하기



```
In [138]: sheet.title = 'D&A practice sheet'
```

```
In [139]: sheet2 = wb.create_sheet('D&A sheet 2')  
wb.save('openpyxl_practice.xlsx')
```

## 시트 이름 변경도 가능

**sheet.title** : 기존의 시트 이름을 변경  
**create\_sheet('값')** : 새로운 시트 이름 지정

## 05 데이터 저장하기

```
In [121]: try:
            wb = openpyxl.load_workbook('D&A.xlsx')
            sheet = wb.active
        except:
            wb = openpyxl.Workbook()
            sheet = wb.active
```

`openxl.load_workbook('파일이름')`  
이미 존재하는 파일을 불러오기  
없으면 에러가 발생

=> 예외처리로 해결

## 과제

네이버 쇼핑에서 원하는 상품을 검색하고, 500개 이상의 데이터를 수집하는 코드를 작성하고, csv파일과 xlsx파일로 저장하시오

- 제품사진, 제품명, 가격, 리뷰수, 구매건수, 찜한수, 사이트 링크를 모두 포함해야함
- 제품사진은 폴더를 따로 만들어 저장한 후 캡처본으로 제출
- 패키지 이용에는 관여하지 않음
- 제일 처음 받는 웹페이지의 주소는 네이버 쇼핑의 첫 화면이어야 함
- 제출 목록 : ipynb파일, csv파일, excel파일, 사진저장 폴더의 캡처본