

# 제 14 강 Web Scraping (2)

# 과업 #2

## 네이버 쇼핑 신라면 가격 최저가 리스트 만들기

카테고리	붕지라면	컵라면
키워드추천	40개	30개 박스 건면
가격	3백원 ~ 7백원 <b>HIT</b>	2천원 ~ 4천원
배송/혜택/색상	무료배송	내일도착 오늘출발

[https://search.shopping.naver.com/search/all?frm=NVSHATC&origQuery=신라면&pagingIndex=1&pagingSize=40&productSet=total&query=신라면&sort=price\\_asc&timestamp=&viewType=list](https://search.shopping.naver.com/search/all?frm=NVSHATC&origQuery=신라면&pagingIndex=1&pagingSize=40&productSet=total&query=신라면&sort=price_asc&timestamp=&viewType=list)

전체	가격비교	네이버페이	백화점/홈쇼핑	쇼핑원도	해외직구
63,416	0	25,037	3,387	13	384

· 네이버 랭킹순 · **낮은 가격순** · 높은 가격순 · 리뷰 많은순 · 리뷰 좋은순 · 등록일순

쇼핑물선택 · 상품다입(전체) · 40개씩 보기

추천 낮은 가격순 으로 정렬되었습니다.(1만원 이상) 적용기준 ON



한봉지 미고랭 핫앤스파이시 라면 400g 야식 신라면 - UnKnown

10,000원

식품 > 라면/면류 > 라면 > 붕지라면

[사는 재미의 발견, 티몬 - 매주 월요일은 티몬데이 / 첫고객 웰컴쿠폰 최대 20% 할인 / 쓸수록 커지는 적립 혜택 슈퍼세이브 / 토스페이 5% 즉시할인]

등록일 2022.04 · 02 찰하기 · 0 신고하기

**TMON** 정보  
클래디넴

배송비 2,500원

구매전보

✓ 네이버 랭킹순 · 낮은 가격순 · 높은 가격순 · 리뷰 많은순 · 리뷰 좋은순 · 등록일순



a.basicList\_link\_\_1MaTN 128.7 × 14

농심 신라면 120g 40개

광고① 26,000원

식품 > 라면/면류 > 라면 > 봉지라면

개당열량 : 500kcal | 무게 : 120g

첫 구매 시 30일 동안 무제한 무료배송 / 로켓배송  
로켓와우 무료배송!

등록일 2021.11. · ♥ 찜하기 149 · 📬 신고하기

```
<div>
  <li class="basicList_item__2XT81 ad">
    <div class="basicList_inner__eYmq">
      <div class="basicList_img_area__a3NRA">...</div>
      <div class="basicList_info_area__17Xyo">
        <div class="basicList_title__3P9Q7">
          <a target="_blank" class="basicList_link__1MaTN" rel="noopener" data-nclick="N=a:1st*A.title,i:2978667155
8,r:1" title="농심 신라면 120g 40개" href="https://adcr.naver.com/adcr?x=K47BYNLRPEi36KCR7yenMf///w==kfaik
Uvcbx...uVXH91KUq/6893bYNLxCSHANADG6szqhCXOyqnxP5wQW391aKPHbsr2/bZAmRJgnwoIVIqX2">농심 신라면 120g 40개</a>
          == $0
        </div>
        <div class="basicList_price_area__1UXXR">...</div>
        <div class="basicList_depth__2QIie">...</div>
        <div class="basicList_desc__2-tko basicList_max__bowiv">...</div>
        <div class="basicList_etc_box__1Jzg6">...</div>
      </div>
      <div class="basicList_mall_area__1IA7R">...</div>
    </div>
  </li>
</div>
```

```
url =  
f'https://search.shopping.naver.com/search/all?frm=NVSHATC&origQuery={keyword}&pagingIndex=1  
&pagingSize=40&productSet=total&query={keyword}&sort=price_asc&timestamp=&viewType=list'  
  
r = requests.get(url, headers=headers)  
r.raise_for_status()  
  
soup = BeautifulSoup(r.text, 'lxml')  
elms = soup.find_all(class_=re.compile(r'^basicList_title'))  
  
for e in elms:  
    title = e.a['title']  
    price = e.next_sibling.find(class_=re.compile('^price_num')).string  
    print(f'{price} : {title}')
```

# 문제점

- JavaScript 를 통한 반응형 웹의 경우, request module 로 제대로 읽어낼 수 없음.

```
▼ <ul class="list_basis">  
  ▼ <div> == $0  
    ▶ <div>...</div>  
    ▶ <div>...</div>  
    ▶ <div>...</div>  
    ▶ <div>...</div>  
    ▶ <div>...</div>  
  </div>  
</ul>
```

# Selenium

- 웹 자동화 테스트 도구.
  - 각종 사용자의 입력(클릭, 텍스트입력, 리스트 선택, ...)을 프로그램으로 처리.
- 인터랙티브 웹 사이트 크롤링에 활용.
- 설치
  - selenium 라이브러리 설치 - `pip install selenium`
  - 자동화를 위한 브라우저 프로그램 설치 - working folder 아래에 복사.
    - Chrome driver - PC 에 설치되어 있는 chrome 의 버전에 맞는 드라이버 설치 필요.
    - <https://chromedriver.chromium.org/downloads>

# Playwright

- Microsoft 사가 개발한 웹 테스트 및 자동화 프레임워크
- Selenium과의 비교 우위
  - 자동 대기 - 엘리먼트들이 준비될 때까지 대기
  - 테스트 레코딩 & 재생
  - 여러 개의 브라우저에서 동시 실행
- 설치 (2단계)
  - `pip install pytest-playwright`
  - `playwright install`

# 네이버 홈

```
from playwright.sync_api import sync_playwright

p = sync_playwright().start()
browser = p.chromium.launch(headless=False)
page = browser.new_page()
page.goto("https://www.naver.com")
print(page.title())
browser.close()
p.stop()
```



```

from playwright.sync_api import sync_playwright
from bs4 import BeautifulSoup
import re
import time

keyword = '신라면'
url =
'https://search.shopping.naver.com/search/all?frm=NVSHATC&origQuery='+keyword+'&pagingIndex=1&pagingSize=10&productSet=total
&query='+keyword+'&sort=price_asc&timestamp=&viewType=list'

p = sync_playwright().start()
browser = p.chromium.launch(headless=False).new_context(
    user_agent='Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/109.0.0.0
Safari/537.36',
    viewport={'width': 1920, 'height': 2048}
)
page = browser.new_page()
page.goto(url)

src_size = 0
while src_size < len(page.content()):
    src_size = len(page.content())
    page.keyboard.press('End')
    time.sleep(1)

soup = BeautifulSoup(page.content(), 'lxml')
elms = soup.find_all(class_=re.compile(r'^basicList_title'))

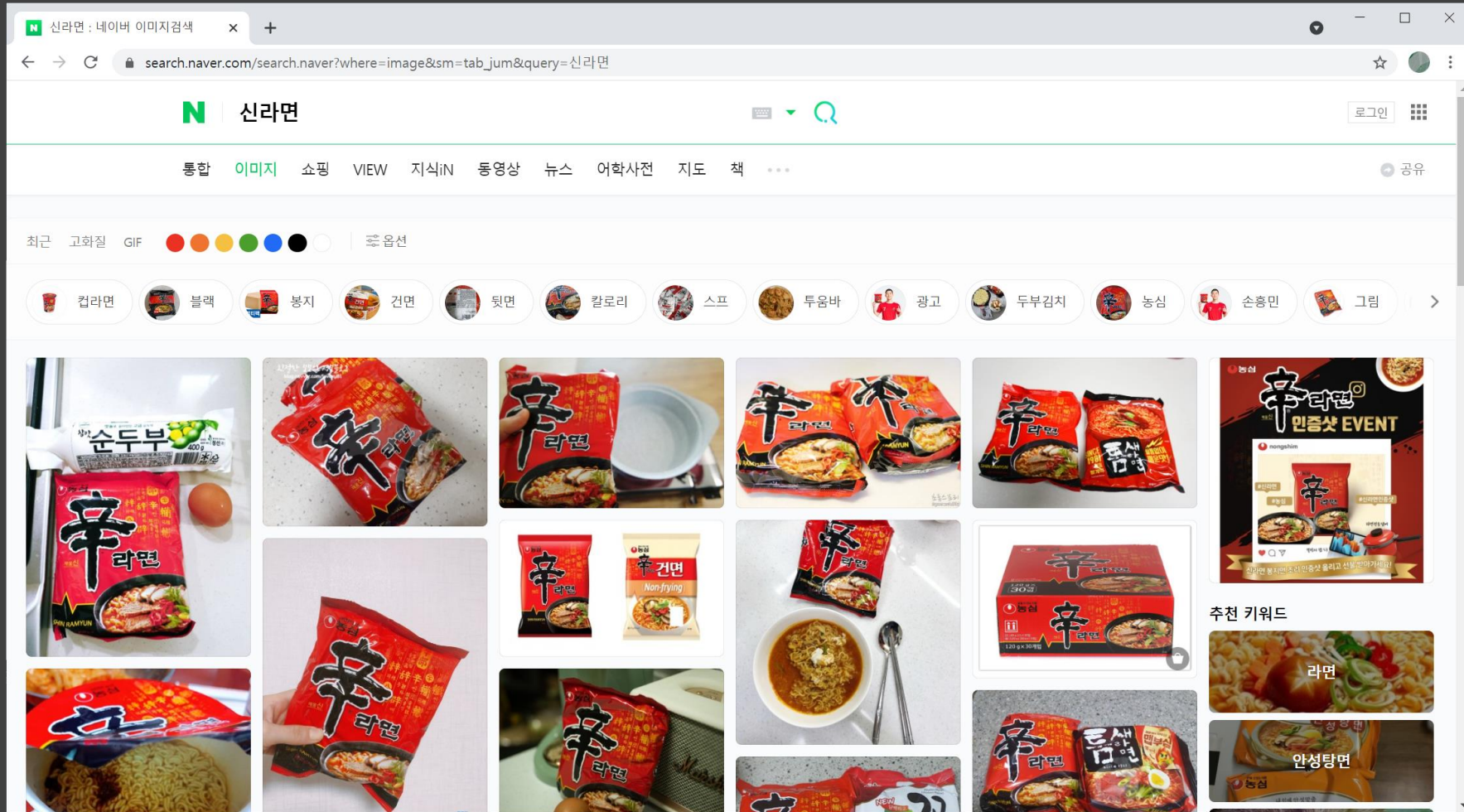
for e in elms:
    title = e.a['title']
    price = e.next_sibling.find(class_=re.compile('^price_num')).text
    print(f'{price=} : {title=}')

browser.close()
p.stop()

```

# 과업 #3.

## ■ 신라면 이미지 자동 다운로드



```
01a4ecce5c96b13d3230451e&g=663960999185414" class="link_thumb _imageBo
x _infoBox" title="[라면추천] 모두가 알고있고 모두가 맛있어하는 신기한 얼
큰라면 신라면 !!!" role="button" aria-pressed="false">
 == $0
<i class="spimg ico_selected"></i>
::after
```

```
keyword = '신라면'
url = 'https://search.naver.com/search.naver?where=image&sm=tab_jum&query=' + keyword

p = sync_playwright().start()
browser = p.chromium.launch(headless=False).new_context(
    user_agent='Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like
    Gecko) Chrome/109.0.0.0 Safari/537.36'
)
page = browser.new_page()
page.goto(url)

soup = BeautifulSoup(page.content(), 'lxml')

elms = soup.select('img._image._listImage')
for n, e in enumerate(elms):
    print(n, e['src'], e)
```

# 로딩 늦어지는 리소스에 대한 처리

```
for n, e in enumerate(elms):  
    caption = e['alt']  
    image_url = e.get('data-lazy-src', e['src'])  
    print(f'{n} {caption} {image_url}')
```

# requests.get 를 이용한 다운로드

```
def download_image(image_url):  
    r = requests.get(image_url)  
    r.raise_for_status()  
    if r.headers['Content-Type'] == 'image/jpeg':  
        fname = f'img_{uuid.uuid4().int}.jpg'  
        with open(fname, 'wb') as wf:  
            print(f'Downloading {fname} .....')  
            wf.write(r.content)  
  
for e in elms:  
    image_url = e.get('data-lazy-src', e['src'])  
    download_image(image_url)
```