



# 분석용 데이터

## 획득 전략

웹 이미지 데이터 확보 기법

# 학습 목표

+ + +

## 학습 목표

- 이미지 데이터의 원본 URL 정보를 수집할 수 있다.
- 인증이 필요 없는 서버에서 이미지 정보와 인증이 필요한 서버에서 이미지 정보를 수집할 수 있다.

## 학습 내용

- 다양한 이미지 데이터 수집하기 - 이론
- 다양한 이미지 데이터 수집하기 - 실습

## 이미지 데이터 수집을 위해 알아야 하는 내용

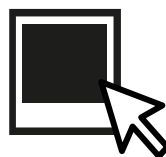
### 1) 작업 순서

01



검색창에 키워드 입력 후 검색

02



이미지 클릭

03



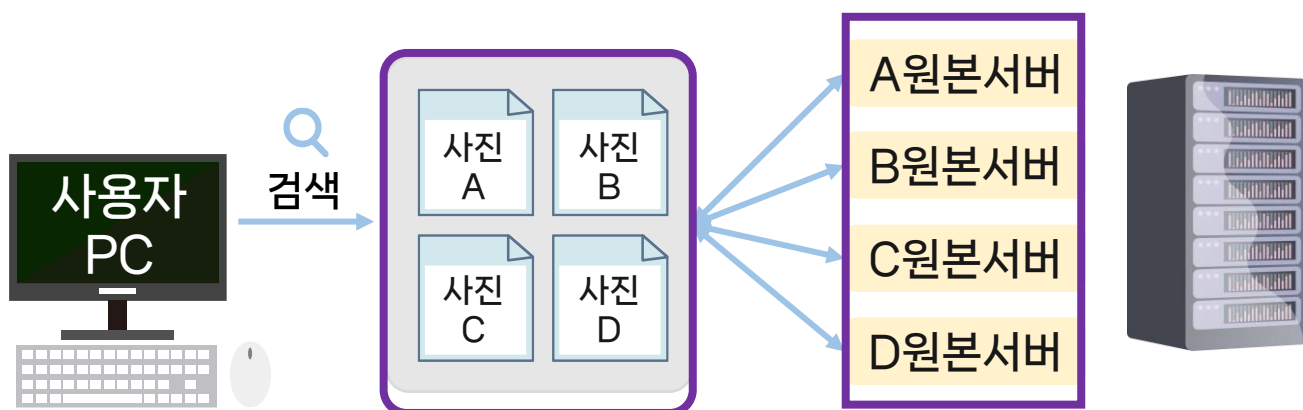
이미지 URL 추출

04



이미지 다운로드 후 저장

### 2) 사진의 원본 URL 주소 추출 → 사진 다운로드



### 이미지 데이터 수집을 위해 알아야 하는 내용

#### 3) 원본 이미지의 URL 주소 필요

 이미지 URL 주소는 img 태그의 src 속성값

```
[<img src='https://cdn.abc.com/aaa.jpg' alt="고양이" .....>]
```

태그이름

URL 주소

#### 4) 주의사항

##### (1) 수집 대상 여부 확인

 수집대상 URL 맞춤

```
[https://search.pstatic.net/common/?src=http%..._8584  
.jpg&type=ofullfill340_600_png]
```

 수집대상아님

```
[data:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///y]
```

#### 주의사항

**수집 대상에 맞는 URL을 구분하여 이미지 데이터 수집 코드를 작성**

### 이미지 데이터 수집을 위해 알아야 하는 내용

#### (2) 이미지 URL에 한글이 들어가면 오류 발생

🔍 오류가 발생하는 경우

[[https://search.pstatic.net/common/fill340\\_파이썬기초.png](https://search.pstatic.net/common/fill340_파이썬기초.png)]

🔍 정상 수집되는 경우

[[https://search.pstatic.net/common/fill340\\_python100.png](https://search.pstatic.net/common/fill340_python100.png)]

#### 주의사항

주소에 **한글**이 포함되어 있을 경우, **변환작업**이 되도록 코드를 작성

#### (3) 웹 브라우저 정보

🔍 웹브라우저 정보가 필요한 서버일 경우

➡ 이미지를 다운로드 받을 때 웹브라우저 정보가 필요할 수 있음

🔍 웹브라우저 정보가 필요한 정보를 주지 않을 경우

➡ HTTPError: HTTPError 403: Forbidden **에러발생**

### 인증정보가 필요 없는 사이트 정보 수집하기

- 한 페이지에 보이는 이미지 개수

```
1 # 다양한 이미지 데이터 수집하기
2 # 필요한 모듈과 라이브러리를 로딩하고 검색어를 입력 받습니다
3 from bs4 import BeautifulSoup
4 from selenium import webdriver
5 from selenium.webdriver.common.by import By
6 from selenium.webdriver.common.keys import Keys
7 from selenium.webdriver.chrome.service import Service
8 import urllib.request
9 import urllib
10 import time
11 import sys
12 import math
13 import os
14 import random
15
16 #사용자에게 필요한 정보들을 입력 받습니다.
17 print("=" * 80)
18 print(" 이 크롤러는 이미지 정보를 수집합니다")
19 print("=" * 80)
20
21 query_txt = input('1. 크롤링할 이미지의 키워드는 무엇입니까?: ')
22 cnt = int(input('2. 크롤링 할 건수는 몇건입니까?: '))
23 real_cnt = math.ceil(cnt / 20)
24 f_dir=input('3. 파일이 저장될 경로만 쓰세요(예: c:\\py_temp\\ ) : ')
25 if f_dir == '' :
26     f_dir = "c:\\py_temp\\"
27
```

### 인증정보가 필요 없는 사이트 정보 수집하기

#### ▪ 새 폴더 생성

```
28 #파일을 저장할 폴더를 생성합니다
29 n = time.localtime()
30 s = '%04d-%02d-%02d-%02d-%02d-%02d' % (n.tm_year, n.tm_mon, n.tm_mday, n.tm_
31
32 img_dir = f_dir+s+'-'+query_txt
33 os.makedirs(img_dir)
34 os.chdir(img_dir)
35
```

#### ▪ 폴더로 이동

```
28 #파일을 저장할 폴더를 생성합니다
29 n = time.localtime()
30 s = '%04d-%02d-%02d-%02d-%02d-%02d' % (n.tm_year, n.tm_mon, n.tm_mday, n.tm_
31
32 img_dir = f_dir+s+'-'+query_txt
33 os.makedirs(img_dir)
34 os.chdir(img_dir)
35
36 #크롬 드라이버를 사용해서 웹 브라우저를 실행합니다.
37 s_time = time.time()
38
39 s = Service("c:/py_temp/chromedriver.exe")
40 driver = webdriver.Chrome(service=s)
41
42 driver.get('https://www.naver.com')
43 e.sleep(random.randrange(2,5))
```

정리

**time.sleep():** 멈추고자 하는 시간(초)만큼 일시정지



## 다양한 이미지 데이터 수집하기 - 실습

### 인증정보가 필요 없는 사이트 정보 수집하기

- 2~5초 사이에 랜덤 시간을 뽑아 일시정지

```
36 #크롬 드라이버를 사용해서 웹 브라우저를 실행합니다.
37 s_time = time.time()
38
39 s = Service("c:/py_temp/chromedriver.exe")
40 driver = webdriver.Chrome(service=s)
41
42 driver.get('https://www.naver.com')
43 time.sleep(random.randrange(2,5))
44
45 element = driver.find_element(By.ID,"query")
46 element.send_keys(query_txt)
47 element.submit()
48
49 #이미지 링크를 선택합니다
50 driver.find_element(By.LINK_TEXT,"이미지").click()
51
52 #스크롤 다운 하스 새서 후 화면 이동하기
53 def scroll_down(driver):
54     driver.execute_script("window.scrollTo(0,document.body.scrollHeight);")
55     time.sleep(3)
```

- 메뉴나 이모티콘 등

```
1 for i in img_src :
2     img_src1= i.find('img')['src']
3     print(img_src1)
4     print()
https://search.pstatic.net/common/?src=http%3A%2F%2Fblogfiles.naver.net%2FmjAyMjExMT1f0SAg%2FMDAxNjY4ODMwMjE1MDk1.PCW99s
RAHy0CvirxEu4qZ4r-LUanhM83Dum77Q9gw48y.fQWWPhFRQncIwCBzPp9it6UHKpLiccvwPdVo4tqaKsg.JPEG.milky_way159%2FIMG_9583.JPG&typ
e=a340
https://search.pstatic.net/common/?src=http%3A%2F%2Fblogfiles.naver.net%2FmjAyMjExMTZf0TEg%2FMDAxNjY4NTkzNzc2MjIz.mJdYae
-cXI3MCGXQNPLhvE8jkSHV0Kv89GvBk0xv7Ksg.ws6DGL04Kwplua4fbTjY_Py69szQAn9q0XYdWxrIL-Ag.JPEG.monami5335%2F77DB5D4A-8225-48D7
4E1F-384A351008AF_1_201_a.jpg&type=a340
data:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAAABAAEAAAIBRAA7
https://search.pstatic.net/common/?src=http%3A%2F%2Fblogfiles.naver.net%2FmjAyMjExMjVfNjAg%2FMDAxNjY5MzQzODM4MjY1.1kepJB
VB1_B4H9iq49LITBIMRWJFLX0n2Iu8tsdbbIEg.fIarIpbqVNpBZ0hQWd7JzzTCsfirqrKsCGHfwI7IdK0g.PNG.birdowitchears39%2F20221125_1127
20.png&type=a340
https://search.pstatic.net/common/?src=http%3A%2F%2Fblogfiles.naver.net%2FmjAyMjEyMDFfMzAw%2FMDAxNjY5ODU5ODc1NjI2.v8DpH1
33H2S2rglTT71L7R1UTE1xiMWHcW3QZ6K9IQg.Znye1ZyhwRyinNsBpZnmBz0w2Zrcwah71SvM21Xx0mMg.PNG.wgeqjripewjr%2F1.PNG&type=a340
data:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAAABAAEAAAIBRAA7
https://search.pstatic.net/common/?src=http%3A%2F%2Fblogfiles.naver.net%2FmjAyMjExMjRfMjg4%2FMDAxNjY5Mjc2NTk3MDk4.-tWDTV
1J0Mo5MuB6t1Qcx63C50K8MTep7JADVa9R8oq.Araskv04RkSi4N4HwwrELRX02Llobxa63tG25KfiiIq.JPEG.themaltese%2F123_%252813%2529.i
```

```
1 # 수집 대상 이미지의 URL 정보 추출
2 for i in img_src :
3     ima_src1= i.find('img')['src']
```



### 인증정보가 필요 없는 사이트 정보 수집하기

#### ▪ 가져올 이미지의 URL 주소

```
10 # 이미지 다운로드하여 저장하기
11 for i in range(0, len(img_src2)) :
12     try :
13         urllib.request.urlretrieve(img_src2[i], str(file_no)+'.jpg')
14     except :
15         continue
16
17     time.sleep(0.5)
18     print("%s 번째 이미지 저장중입니다=====" % file_no)
19
20     file_no += 1
21
22     if file_no > cnt :
23         break
```

#### ▪ 저장할 파일 이름

```
10 # 이미지 다운로드하여 저장하기
11 for i in range(0, len(img_src2)) :
12     try :
13         urllib.request.urlretrieve(img_src2[i], str(file_no)+'.jpg')
14     except :
15         continue
16
17     time.sleep(0.5)
18     print("%s 번째 이미지 저장중입니다=====" % file_no)
19
20     file_no += 1
21
22     if file_no > cnt :
23         break
```

### 인증정보가 필요 없는 사이트 정보 수집하기

#### ▪ 인코딩 변환

```
1 # 이미지 이름에 한글이 들어갈 경우 조치방법
2 for i in range(0, len(img_src2)) :
3     file_no += 1
4     urllib.request.urlretrieve(urllib.parse.quote(img_src2[i].encode('utf8'), '/:'), \
5                               str(file_no)+'.jpg')
6
7     time.sleep(0.5)
8     print("%s 번째 이미지 저장중입니다=====" % file_no)
9
10    file_no += 1
11
12    if file_no > cnt :
13        break
```

```
1 # 요약 정보를 출력합니다
2 e_time = time.time()
3 t_time = e_time - s_time
4
5 store_cnt = file_no - 1
6
7 print("-" * 70)
8 print("총 소요시간은 %s 초 입니다 " % round(t_time, 1))
9 print("총 저장 건수는 %s 건 입니다 " % store_cnt)
```

## 인증정보가 필요한 사이트 정보 수집하기

```
4 from bs4 import BeautifulSoup
5 from selenium import webdriver
6 from selenium.webdriver.common.by import By
7 from selenium.webdriver.common.keys import Keys
8 from selenium.webdriver.chrome.service import Service
9 import urllib.request
10 import urllib
11 import time
12 import math
13 import os
14 import random
15
16 #Step 2. 필요한 정보를 입력 받습니다.
17 print("=" * 80)
18 print(" pixabay 사이트에서 이미지를 검색하여 수집하는 크롤러 입니다 ")
19 print("=" * 80)
20
21 query_txt = input('1. 크롤링할 이미지의 키워드는 무엇입니까?: ')
22 cnt = int(input('2. 크롤링 할 건수는 몇건입니까?: '))
23 real_cnt = math.ceil(cnt / 100) # 실제 크롤링 할 페이지 수
24 f_dir=input('3.파일이 저장될 경로만 쓰세요(예: c:\\py_temp\\ ) : ')
25 if f_dir == '' :
26     f_dir = "c:\\py_temp\\"
27
28 print("\n")
29 print("요청하신 데이터를 수집 중이오니 잠시만 기다려 주세요~~^^")
30
31 #Step 3. 파일을 저장할 폴더를 생성합니다
32 n = time.localtime()
33 s = '%04d-%02d-%02d-%02d-%02d-%02d' % (n.tm_year, n.tm_mon, n.tm_mday, n.tm_hour, n.tm_min,
34
주의사항 dir+s+'-'+query_txt
      (img_dir)
```

무료 이미지라도 사용 용도에 따라 **저작권 위반**이 될 수 있음