



Python News Crawling TIL 2

2022-03-08

[2탄] 쉽게 따라하는 네이버 뉴스 크롤링(python) - title, URL 가져오기

"본 포스팅은 네이버 웹 크롤링 실제 python 코드를 작성하는 2탄입니다. 전 단계인 수행계획을 확인하고 싶으신 분들은 아래링크(1탄)를 참고해주세요 :)" 네이버 웹 페이지 구성이 바뀌어 내용, 코드 수정(2021...

https://everyday-tech.tistory.com/3

input	1. 검색어 ex) "코로나" 2. 필요한 뉴스 기사 수 ex) 25										
사용 툴	Python(requests, BeautifulSoup 패키지) - Requests : Request를 보내 웹 페이지 소스를 받음 - BeautifulSoup : 웹 페이지의 소스를 parsing하여 원하는 정보를 찾음										
	<table border="1"> <thead> <tr> <th>사용 툴</th><th>HTML 상 위치</th></tr> </thead> <tbody> <tr> <td>1. 홈페이지의 뉴스 기사들이 모여있는 태그</td><td>body 태그 & class속성값="list_news"</td></tr> <tr> <td>2. 하나의 뉴스 기사에 대한 태그</td><td>li 태그 & id속성값="sp_news"으로 시작</td></tr> <tr> <td>3. 2번 태그의 본문, title, URL이 들어 있는 태그</td><td>div 태그 & class속성값="news_item"</td></tr> <tr> <td>4. 3번 태그의 뉴스 제목, URL이 들어 있는 태그</td><td>a 태그 & class속성값="news_title"</td></tr> </tbody> </table>	사용 툴	HTML 상 위치	1. 홈페이지의 뉴스 기사들이 모여있는 태그	body 태그 & class속성값="list_news"	2. 하나의 뉴스 기사에 대한 태그	li 태그 & id속성값="sp_news"으로 시작	3. 2번 태그의 본문, title, URL이 들어 있는 태그	div 태그 & class속성값="news_item"	4. 3번 태그의 뉴스 제목, URL이 들어 있는 태그	a 태그 & class속성값="news_title"
사용 툴	HTML 상 위치										
1. 홈페이지의 뉴스 기사들이 모여있는 태그	body 태그 & class속성값="list_news"										
2. 하나의 뉴스 기사에 대한 태그	li 태그 & id속성값="sp_news"으로 시작										
3. 2번 태그의 본문, title, URL이 들어 있는 태그	div 태그 & class속성값="news_item"										
4. 3번 태그의 뉴스 제목, URL이 들어 있는 태그	a 태그 & class속성값="news_title"										

수행 단계

STEP 1. 소스 조사 : 제공 사이트 조사, 제공 정보 조사, 확보 가능 정보 확인

STEP 2. 웹 구성 체계 확인 : HTML 구조 확인, 제공 정보 확인, 크롤링 가능 여부 확인

STEP 3. 크롤링 진행 : parsing 방안 및 위치 확인, request 방법 확인, data 저장 형태 설계

STEP 4. 최종 데이터 생성 : 데이터 저장 형태, 데이터 저장

소스 조사 및 웹 구성 체계 확인 결과

- 네이버 뉴스의 기사 URL 형태 (우크라이나를 키워드로 검색 시)
 - https://search.naver.com/search.naver?where=news&ie=utf8&sm=nws_hy&query=우크라이나
- 네이버 뉴스 HTML 구성 (뉴스 title, URL 찾을 시)
 - source : https://search.naver.com/search.naver?where=news&ie=utf8&sm=nws_hy&query="키워드"
 - input : 검색어 (예 : 우크라이나), 필요한 뉴스 기사 수 (예 : 25)
 - output : 뉴스 제목, URL 이 담긴 엑셀
- 크롤링 사용 툴 : Python (Requests, BeautifulSoup 패키지)
 - Requests : request 를 보내 웹 페이지 소스 받음
 - Beautifulsoup : 웹 페이지의 소스를 parsing 하여 원하는 정보를 찾음

STEP 3. 크롤링 진행

- 필요 환경 및 패키지 설치
 - python

```
C:\Users\multicampus>python
Python 3.7.9(tags/v3.7.9:13c94747c7, A
Type "help", "copyright", "credits" or
>>>
```

- requests, bs4, re, pandas 패키지 사용
 - requests : 웹 페이지 소스 추출 (HTML)

- bs4 : HTML 파싱, 필요 태그 및 소스 추출
- re : 조건부 문자열 (정규 표현식), 태그 탐색 시 일반화 조건을 사용하기 위함
- pandas : 데이터 프레임, 엑셀 변환
- 코드 작성
 - 패키지 importing

```
# 패키지 importing

import requests
from pandas import DataFrame
from bs4 import BeautifulSoup
import re
from datetime import datetime
import os
```

```
# 편집 모드 단축키
shift + enter : run cell, select below
ctrl + enter : run cell
alt + enter : run cell, insert below

# mode 변경
ctrl + m / esc

# Command 모드 단축키
b : 셀 추가 (아래로 추가)
a : 셀 추가 (위로 추가)
dd : 해당 셀 삭제
```

- 현재 시간 저장

```
# 현재 시간 저장 : 나중에 output 으로 엑셀 저장 시 크롤링 한 날짜, 시간을 파일명에 넣기 위해 저장하는 변수

date = str(datetime.now())
date = date[:date.rfind(':')].replace(' ', '_')
date = date.replace(':', '시') + '분'
```

- input 생성

```
query = input('검색 키워드를 입력하세요 : ')
query = query.replace(' ', '+')
news_num = int(input('총 필요한 뉴스기사 수를 입력해주세요(숫자만 입력) : '))
```

- 요청할 URL 생성 및 요청

```
# 요청할 URL 생성 및 요청 :
# query 로 받은 키워드를 URL 조건절 중 키워드에 해당하는 변수에 대응시켜 요청 URL 을 만든다
# 그리고 requests 패키지의 get함수를 이용하여 HTML 코드를 받아온다
# 받은 코드를 bs4 BeautifulSoup 함수를 이용하여 파싱한다

news_url = 'https://search.naver.com/search.naver?where=news&sm=nws_hyq&query={}'
req = requests.get(news_url.format(query))
soup = BeautifulSoup(req.text, 'html.parser')
```

- 원하는 정보를 담은 변수 생성 (딕셔너리)

```
# 원하는 정보를 담은 변수 생성 (딕셔너리) :
# 뉴스 기사 정보를 저장할 딕셔너리를 생성 (key : 번호, value : 뉴스 기사 정보)
# idx : 현재 뉴스의 번호
# cur_page : 네이버 뉴스의 웹 페이지, 추출하려는 기사 수가 현재 페이지에 있는 기사보다 많은 경우 다음페이지로 넘어가야 하기 때문에 현 페이지 번호를 기억하5

news_dict = {}
idx = 0
cur_page = 1
```

- parsing 한 HTML 코드에서 원하는 정보 탐색 (뉴스 기사 title, URL)


```

xlsx_file_name = '네이버뉴스_{}_{}.xlsx'.format(query, date)
news_df.to_excel(xlsx_file_name)

# 저장을 완료한 폴더를 띄움
print('엑셀 저장 완료 | 경로 : {}'.format(folder_path, xlsx_file_name))
os.startfile(folder_path)

```

STEP 4. 최종 데이터 생성

```

area_list = [m.find('div', {'class' : 'news_area'}) for m in m_list] # 각 뉴스
a_list = [area.find('a', {'class' : 'news_tit'}) for area in area_list] # 각 뉴스

for n in a_list[:min(len(a_list), news_num-idx)]:
    news_dict[idx] = {'title' : n.get('title'), 'url' : n.get('href')}
    idx += 1
    cur_page += 1

# 현재 수집한 뉴스 기사 수가 부족한 경우 다음 페이지로 넘어가야 하므로 다음 페이지
# 하위에 존재하는 a 태그 내부에 페이지 번호와 URL 정보 존재
# cur_page 변수와 일치하는 페이지 번호의 URL 가져옴
pages = soup.find('div', {'class' : 'sc_page_inner'})
next_page_url = [p for p in pages.find_all('a') if p.text == str(cur_page)][0].get('href')
req = requests.get('https://search.naver.com/search.naver' + next_page_url)
soup = BeautifulSoup(req.text, 'html.parser')

```

크롤링 중...

```

In [51]: # 데이터 프레임 변환 및 저장

print('크롤링 완료')

# 크롤링한 뉴스 정보가 담긴 딕셔너리를 데이터 프레임으로 변환
print('데이터프레임 변환')
news_df = DataFrame(news_dict).T

# 크롤링한 키워드와 크롤링 날짜를 엑셀 파일 명으로 하여 저장
folder_path = os.getcwd()
xlsx_file_name = '네이버뉴스_{}_{}.xlsx'.format(query, date)
news_df.to_excel(xlsx_file_name)

# 저장을 완료한 폴더를 띄움
print('엑셀 저장 완료 | 경로 : {}'.format(folder_path, xlsx_file_name))
os.startfile(folder_path)

크롤링 완료
데이터프레임 변환
엑셀 저장 완료 | 경로 : C:\Users\multicampus\Desktop\GITHUB\boonmoja\crawling\네이버뉴스_우크라이나_2022-03-08_14시28분.xlsx

```

