

우리에게 보여지는 네이버 블로그 URL은 사실 가짜 주소인 것을 알고 계신가요? 이번 포스팅에서는 네이버 블로그의 도서 리뷰를 크롤링하는 과정에서 학습한 내용을 공유합니다.

## 네이버 추천도서 정보 웹 크롤링

추천모델을 생성하는 연습용 프로젝트를 진행하고 있습니다. 네이버 책의 소설 > 나라별 소설 카테고리의 추천도서 969권의 정보를 수집, Word2Vec 모델 기반으로 추천 모델을 생성하고자 합니다.

수집하고자 하는 책 정보의 항목은 다음과 같습니다.

- 책 제목 텍스트
- 책 상세 웹페이지 URL 텍스트
- 상세 웹페이지의 '책소개' 항목 텍스트
- 상세 웹페이지의 '책 속으로' 항목 텍스트
- 상세 웹페이지의 '네티즌 리뷰' 항목 텍스트
  - 개인들의 네이버 블로그에 작성된 책 리뷰
  - 리뷰 총합이 100건을 초과할 경우 100건만 수집

문제는 리뷰 텍스트의 수집에서 발생했습니다. Selenium(chromedriver)와 BeautifulSoup 등 다양한 방법으로 접근해 봐도, 본문 텍스트를 가져오지 못했는데요.

네이버 API를 사용하여 urllib.request.Request 객체에 API 사용자 인증 정보를 header 추가하여도 마찬가지였습니다. 무엇이 문제였을까요?

# 가짜 URL, 진짜 URL

네이버 블로그에 업로드된 포스트는 각각 가짜 주소와 진짜 주소를 가지고 있습니다. 가짜 주소는 우리가 일반적으로 네이버 블로그 포스트에 접근할 때 보여지는 주소이며, 웹 크롤링이 불가능하도록 본문 등 콘텐츠 텍스트가 숨겨져 있습니다.

도저히 콘텐츠 텍스트를 얻을 수가 없었던 이유입니다. 가짜 URL에서 정보를 가져오려고 시도했던 것이니까요. 이와 반대로, 진짜 URL은 크롤링이 가능한 주소이며, 숨겨져 있습니다.

어디에 숨겨져 있을까요? 아래 도서 리뷰 포스트를 예시로 들어보겠습니다.

- <https://blog.naver.com/box1042/220907410772>

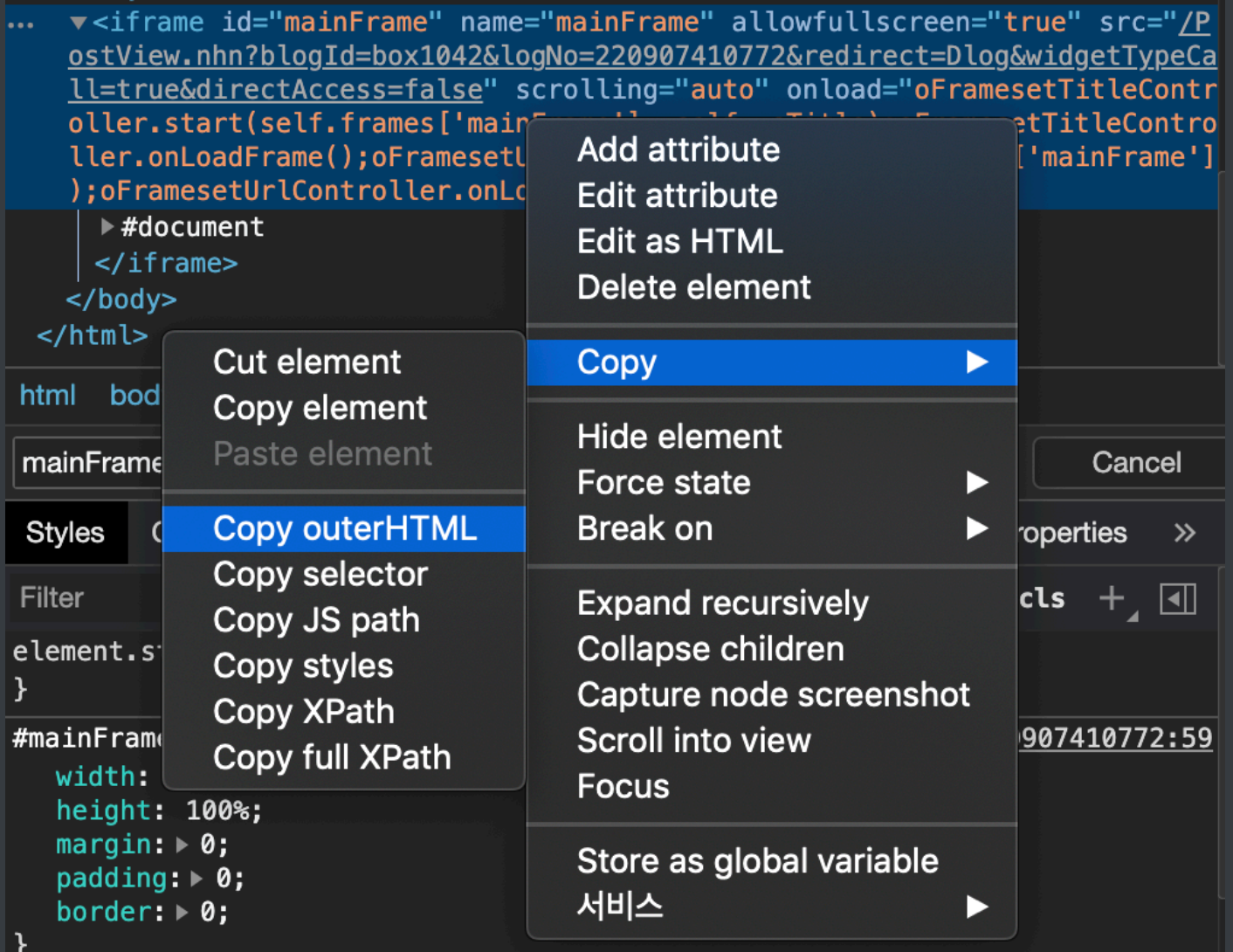
크롬 브라우저로 위 주소에 접속하여 마우스 우클릭 > 검사(investigation)를 찍어 보면 해당 페이지의 코드를 확인할 수 있습니다. Ctrl+F로 'mainFrame'을 검색하면 다음과 같은 결과를 확인할 수 있는데요.

```
... ▾ <body> == $0
  ▾ <iframe id="mainFrame" name="mainFrame" allowfullscreen="true" src="/PostView.nhn?blogId=box1042&logNo=220907410772&redirect=Dlog&widgetTypeCall=true&directAccess=false" scrolling="auto" onload="oFramesetTitleController.start(self.frames['mainFrame'], self, sTitle);oFramesetTitleController.onLoadFrame();oFramesetUrlController.start(self.frames['mainFrame']);oFramesetUrlController.onLoadFrame();">
    ▶ #document
    </iframe>
  </body>
</html>
```

html body

mainFrame 2 of 2 ^ v Cancel

위 코드를 마우스 우클릭하여 Copy > Copy outerHTML하면, id와 name 속성을 mainFrame으로 가지는 iframe에 대한 코드를 확인할 수 있습니다.



```
<iframe id="mainFrame" name="mainFrame" allowfullscreen="true"
src="/PostView.nhn?
blogId=box1042&logNo=220907410772&from=search&redirect=Log&a
mp;widgetTypeCall=true&topReferer=https%3A%2F%2Fsearch.naver.com%2Fs
earch.naver%3Fquery%3D%255B%25EC%25A0%2595%25EC%2583%2581%25EC%2596%25B4
%25ED%2595%2599%25EC%259B%2590%2520%25EC%259D%25B8%25EC%25B2%259C%25EB%2
582%25A8%25EB%258F%2599%25EB%25B6%2584%25EC%259B%2590%28%25EA%25B5%25AC%
25EC%259B%2594%25EB%258F%2599%29%255D%25EC%2584%259C%25EC%259A%25B8%25EB
%258C%2580%2520%25EA%25B6%258C%25EC%259E%25A5%2520%25EB%258F%2584%25EC%2
584%259C%2520-
%2520%25EC%2595%2588%25EB%2582%2598%2520%25EC%25B9%25B4%25EB%25A0%2588%2
5EB%25A6%25AC%25EB%2582%2598%2522%26nso%3D%26where%3Dblog%26sm%3Dtab_viw
.all&directAccess=false" scrolling="auto"
onload="oFramesetTitleController.start(self.frames['mainFrame'], self,
sTitle);oFramesetTitleController.onLoadFrame();oFramesetUrlController.st
art(self.frames['mainFrame']);oFramesetUrlController.onLoadFrame()">
</iframe>
```

위 코드에서 src라는 속성에 부여된 텍스트가 있습니다. "/PostView.nhn?blogId=..."로 시작하는 텍스트인데요. 이 텍스트가 바로 진짜 주소값의 일부를 가리킵니다. 정확히 말하면, <https://blog.naver.com> 다음에 오는 주소값입니다.

정리하면, <https://blog.naver.com/box1042/220907410772> 라는 네이버 블로그 포스트의 웹크롤링 가능한 진짜 주소는 다음과 같게 됩니다.

```
https://blog.naver.com/PostView.nhn?
blogId=box1042&logNo=220907410772&from=search&redirect=Log&
mp;widgetTypeCall=true&topReferer=https%3A%2F%2Fsearch.naver.com%2Fs
earch.naver%3Fquery%3D%255B%25EC%25A0%2595%25EC%2583%2581%25EC%2596%25B4
%25ED%2595%2599%25EC%259B%2590%2520%25EC%259D%25B8%25EC%25B2%259C%25EB%2
582%25A8%25EB%258F%2599%25EB%25B6%2584%25EC%259B%2590%28%25EA%25B5%25AC%
25EC%259B%2594%25EB%258F%2599%29%255D%25EC%2584%259C%25EC%259A%25B8%25EB
%258C%2580%2520%25EA%25B6%258C%25EC%259E%25A5%2520%25EB%258F%2584%25EC%2
584%259C%2520-
%2520%25EC%2595%2588%25EB%2582%2598%2520%25EC%25B9%25B4%25EB%25A0%2588%2
5EB%25A6%25AC%25EB%2582%2598%2522%26nso%3D%26where%3Dblog%26sm%3Dtab_viw
.all&directAccess=false
```

이제 가짜 주소로부터 진짜 주소를 획득하는 방법을 알았습니다. 아래와 같은 코드를 작성하여, 네이버 블로그 포스트의 본문 내용을 가져와 봅시다.

```
from bs4 import BeautifulSoup
import requests
import urllib.request
import urllib.error
import urllib.parse
import re
import json

naver_client_id = "#####"
naver_client_secret = "#####"
```

```
def get_naver_blog_text(client_id_, client_secret_, keyword):
    """
    client_id_ : 네이버 검색 API ID
    client_secret_ : 네이버 검색 API password
    keyword : 검색어 (네이버 블로그 포스트 제목)
    """
```

```

# naver API info
client_id = client_id_
client_secret = client_secret_

encText = urllib.parse.quote(keyword)

# 검색 쿼리를 전송할 url address
url = "https://openapi.naver.com/v1/search/blog?query=" + encText #
json 결과
# url = "https://openapi.naver.com/v1/search/blog.xml?query=" +
encText # xml 결과

# 쿼리 url 전송
request = urllib.request.Request(url)
# Naver API 인증
request.add_header("X-Naver-Client-Id",client_id)
request.add_header("X-Naver-Client-Secret",client_secret)

# 서버 응답 수신
response = urllib.request.urlopen(request)
rescode = response.getcode()

# 서버 응답 코드 200 (정상)
if(rescode==200):
    # 서버가 전송한 내용으로부터 블로그 포스트 url 추출
    response_body = response.read()
    response_body.decode('utf-8')
    response_body_dict = json.loads(response_body.decode('utf-8'))

    # 블로그 포스트 제목으로 검색한 결과의 첫번째 아이템의 url
    blog_post_url = response_body_dict['items'][0]
    ['link'].replace("&", "")

    # 추출한 블로그 포스트 url로부터 BeautifulSoup 생성
    get_blog_post_content_code = requests.get(blog_post_url)

```

```
get_blog_post_content_text = get_blog_post_content_code.text
get_blog_post_content_soup =
BeautifulSoup(get_blog_post_content_text, 'lxml')

# BeautifulSoup으로부터 웹크롤링 가능한 진짜 url 획득
real_blog_post_url = "http://blog.naver.com" +
get_blog_post_content_soup.select('#mainFrame')[0].get('src')

# 진짜 url로부터 BeautifulSoup 생성
get_real_blog_post_content_code =
requests.get(real_blog_post_url)
get_real_blog_post_content_text =
get_real_blog_post_content_code.text
get_real_blog_post_content_soup =
BeautifulSoup(get_real_blog_post_content_text, 'lxml')

# 포스트의 본문 텍스트를 저장할 리스트
contents = []

# 진짜 url 웹페이지의 본문 태그 내용에 대하여 for loop
for content in
get_real_blog_post_content_soup.select('div#postViewArea'):
    # 본문 텍스트 추출
    blog_post_content_text = content.get_text()

    # HTML 태그 제거
    remove_html_tag = re.compile('<.*?>')
    blog_post_full_contents = str(blog_post_content_text)

    # 리스트에 HTML 태그 제거한 본문 텍스트 저장
    contents.append(blog_post_full_contents)

    # 리스트에 존재하는 텍스트를 string으로 합쳐서 반환
    return ' '.join(contents)

# 서버 응답 코드가 200이 아님 (비정상)
```



```
# 에러 코드를 출력하고 None 반환
print("Error Code:" + rescode)
return None
```

```
get_naver_blog_text(naver_client_id, naver_client_secret,  
                    "[정상어학원 인천남동분원(구월동)]서울대 권장 도서 - 안나 카레리나")
```

[illegible]

- <https://blog.naver.com/here-now-lee/222215012401>



동일한 방식으로 위에서 작성한 함수를 실행했지만 None이 반환되었습니다. 위 블로그 포스트의 진짜 URL을 찾아서, 브라우저를 통해 직접 접근을 시도해 보았습니다. 그 결과, 다음과 같은 팝업을 확인했습니다.

**blog.naver.com 내용:**

삭제되었거나 존재하지 않는 게시물입니다.

확인

이 경우, 아래 2가지에 대하여 실패했다고 생각할 수 있습니다.

1. 네이버 검색 API를 통해 접근했으나 실패
2. 진짜 URL 주소로 직접 접근했으나 실패

이러한 실패 경우에 대해서도 크롤링을 수행할 수 있는 방법은 추가적인 고민이 필요해 보입니다. 이상으로 오늘의 포스팅을 마칩니다.

## Reference

- [네이버 검색 API 사용하기](#)
- [네이버 블로그의 특정 키워드 크롤링하기](#)
- [파이선 크롤링 체험기](#)
- [jy617lee/naver\\_blog\\_crawler](#)