

```

1 Lab. Using BeautifulSoup
2
3 1. BeautifulSoup module 이용하기
4 1) pip install BeautifulSoup4
5 2) Scraping을 위한 module
6 3) 루이스 캐럴의 [이상한 나라의 앨리스]에 나오는 동명의 시에서 이름을 따왔다.
7 4) 일반적으로 HTML tag들이 start tag와 end tag가 서로 pair 되지 않을 경우가 많다.
8 5) pair 되지 않아도 아름답게 처리해 주는 module
9 6) https://www.crummy.com/software/BeautifulSoup/bs4/doc/
10
11 from bs4 import BeautifulSoup
12 import requests
13
14 html_data = requests.get('https://www.naver.com')
15 soup = BeautifulSoup(html_data.text, 'html.parser')
16 soup.title
17 -----
18 <title>NAVER</title>
19
20 7) 첫번째 인자는 Python 객체로 바꿀 string을 넣어주고 두 번째 인자로 parser를 넣는다.
21 8) Parser란, 원시 코드인 순수 문자열 객체를 해석할 수 있도록 분석하는 것을 의미한다.
22 9) Python에서 사용되는 parser는 다음과 같다.
23 a. lxml
24 b. html5lib
25 c. html.parser
26
27 10) 각 parser의 장단점을 알아보자.
28 a. lxml
29 -XML 해석이 가능한 parser이다.
30 -Python 2.x와 3.x 모두 지원 가능하다.
31 -다른 parser에 비해 매우 빠른 속도로 처리한다.
32 -그 이유는 C언어로 구현되어 있기 때문이다.
33 b. html5lib
34 -웹 브라우저 방식으로 HTML을 해석한다.
35 -하지만 처리 속도가 매우 느리다는 단점이 있다.
36 -그리고 2.x 전용이다.
37 c. html.parser
38 -최신 버전의 Python에서는 사용이 불가
39
40
41 2. 각 parser 비교하기
42 1) lxml
43
44 from bs4 import BeautifulSoup
45 html = """<p>test</p>"""
46 soup = BeautifulSoup(html, 'lxml')
47 print(soup)
48 -----
49 <html><body><p>test</p></body></html>
50
51
52 from bs4 import BeautifulSoup
53 html = """<html><p>test</p></html>"""
54 soup = BeautifulSoup(html, 'lxml')
55 print(soup)
56
57 html = """<body><p>test</p></body>"""
58 soup = BeautifulSoup(html, 'lxml')
59 print(soup)
60 -----
61 <html><body><p>test</p></body></html>
62 <html><body><p>test</p></body></html>
63
64 -html과 body tag가 포함된 형태로 만들어 준다.
65
66
67 2) html5lib
68 $ pip install html5lib
69
70 from bs4 import BeautifulSoup
71 html = """<p>test</p>"""
72 soup = BeautifulSoup(html, 'html5lib')
73 print(soup)
74 -----
75 <html><head></head><body><p>test</p></body></html>
76
77 -html5lib도 html처럼 해석하기 때문에 html, head, body tag가 포함된 형태로 만들어 준다.
78
79
80 3. 한빛미디어 책 제목 읽어오기
81
82 hanbit = requests.get('http://www.hanbit.co.kr/media/')
83 soup = BeautifulSoup(hanbit.text, 'html.parser')
84 soup

```

```

85 -----
86 <!DOCTYPE html>
87 <html lang="ko">
88 <head>
89 <!--[if lte IE 8]>
90 <script>
91     location.replace('/support/explorer_upgrade.html');
92 </script>
93 <![endif]-->
94 <!-- Google Tag Manager -->
95 <script>(function(w,d,s,l,i){w[l]=w[l]||[];w[l].push({'gtm.start':
96     new Date().getTime(),event:'gtm.js'});var f=d.getElementsByTagName(s)[0],
97     j=d.createElement(s),dl=!!'dataLayer'?'&l='+l:'';j.async=true;j.src=
98     'https://www.googletagmanager.com/gtm.js?id='+i+dl;f.parentNode.insertBefore(j,f);
99     })(window,document,'script','dataLayer','GTM-W9D5PM3');

```

```

104 for book in soup.find_all('p', class_='book_tit'):
105     print(book.find('a').text)
106 -----

```

```

107 마이크로서비스 아키텍처 구축 가이드
108 맛있는 디자인 프리미어 프로&애프터 이펙트 CC 2023
109 맛있는 디자인 포토샵&일러스트레이터 CC 2023
110 실전에서 바로 쓰는 Next.js
111 (NO.1 영상 편집 유튜버) 비도클래스의 유튜브 영상 편집 with 프리미어 프로
112 이것이 Windows Server다(개정판)
113 업무에 바로 쓰는 AWS 입문
114 성공하는 유튜브 채널은 따로 있다
115 프로젝트 매니저는 무슨 일을 하고 있을까
116 우리가 사랑한 한국 PC 게임
117 혼자 공부하는파이썬(개정판)
118 혼자 공부하는 머신러닝+딥러닝
119 회사에서 바로 통하는 실무 엑셀+파워포인트+워드&한글(모든 버전 사용 가능, 개정판)
120 혼자 공부하는 컴퓨터 구조+운영체제
121 혼자 공부하는 C 언어
122 맛있는 디자인 포토샵&일러스트레이터 CC 2022
123 혼자 해도 프로 작가처럼 잘 그리는 아이패드 드로잉 with 프로크리에이트
124 구글 엔지니어는 이렇게 일한다
125 국내 최대 이모티콘 승인 작가 씨엠제이가 알려주는 승인율 99.9% 이모티콘 만들기
126 이것이 취업을 위한 코딩 테스트다 with 파이썬
127
128
129

```

4. Naver 영화 평점 Scraping 하기

```

130
131
132 from bs4 import BeautifulSoup
133
134 html_data = requests.get('https://movie.naver.com/movie/point/af/list.nhn?page=1')
135 soup = BeautifulSoup(html_data.text, 'html.parser')
136 titles = soup.find_all(class_='movie')
137
138 title_list = []
139 for title in titles:
140     print(title.text)
141 -----
142 메간
143 애프터썸
144 터미네이터
145 열한번째 엄마
146 접속
147 영웅
148 바빌론
149 가재가 노래하는 곳
150 더 퍼스트 슬램덩크
151 더 퍼스트 슬램덩크
152
153 title_list = []
154 for title in titles:
155     title_list.append(title.text)
156
157 point_list = []
158 points = soup.find_all(class_='list_netizen_score')
159 for point in points:
160     point_list.append(point.find('em').text)
161
162 review_list = []
163 reviews = soup.find_all(class_='title')
164 for review in reviews:
165     rev = review.text
166     rev = rev.strip()
167     rev = rev.replace('\t', '')
168     rev = rev.replace('\n', '')

```

```

169     rev = rev.replace('신고', '')
170     review_list.append(rev)
171
172 df = pd.DataFrame(title_list, columns=['Title'])
173 df['Point'] = point_list
174 df['Review'] = review_list
175 df
176 -----
177      Title      Point      Review
178 0 자백  9 자백별점 - 총 10점 중9기대이상입니다. 반전이 일단 최고네요. 간만에 긴장감있게...
179 1 타이타닉 10 타이타닉별점 - 총 10점 중10여운의 끝을 보여주는 영화....
180 2 자백  7 자백별점 - 총 10점 중7생각보다 재밌게 봤습니다. 기대 안했는데 영상미가 돋보였...
181 3 메간  8 메간별점 - 총 10점 중8근래 본 공포영화 중에선 가장 불만 했습니다.
182 4 누구의 딸도 아닌 해원 10 누구의 딸도 아닌 해원별점 - 총 10점 중10찌질하다는 것은 그릇이 작은 것이다....
183 5 압구정 8 압구정별점 - 총 10점 중8다들 재미없다고 하시는데.. 저는 그냥 아무생각없이 웃...
184 6 영웅 10 영웅별점 - 총 10점 중10이거 극장가서 안보면 완전 후회합니다 초강추요영상미 넘...
185 7 더 퍼스트 슬램덩크 10 더 퍼스트 슬램덩크별점 - 총 10점 중10감동 그 자체입니다..
186 8 타이타닉 10 타이타닉별점 - 총 10점 중10ost 흘러 나올 때마다 울었다..
187 9 타이타닉 10 타이타닉별점 - 총 10점 중10영화관에서는 처음 보네요. 재밌고 감동적이에요
188
189
190

```

5. Naver 평점 1page부터 100page까지 scraping 하기

```

192
193 from bs4 import BeautifulSoup
194 import requests
195 import pandas as pd
196
197 url = 'https://movie.naver.com/movie/point/af/list.nhn?page='
198
199 title_list = []
200 point_list = []
201 review_list = []
202
203 for pge in range(1, 101):
204     url = url + str(pge)
205     print(url)
206     html_data = requests.get(url)
207     soup = BeautifulSoup(html_data.text, 'html.parser')
208     titles = soup.find_all(class_='movie')
209     points = soup.find_all(class_='list_netizen_score')
210     reviews = soup.find_all(class_='title')
211     for title in titles:
212         title_list.append(title.text)
213     for point in points:
214         point_list.append(point.text)
215     for review in reviews:
216         rev = review.text
217         rev = rev.strip()
218         rev = rev.replace('\t', '')
219         rev = rev.replace('\n', '')
220         rev = rev.replace('신고', '')
221         review_list.append(rev)
222     url = url.split('=')[0] + '='
223
224 df = pd.DataFrame(title_list, columns=['Title'])
225 df['Point'] = point_list
226 df['Review'] = review_list
227
228 df.info()
229 -----
230 https://movie.naver.com/movie/point/af/list.nhn?page=1
231 https://movie.naver.com/movie/point/af/list.nhn?page=2
232 https://movie.naver.com/movie/point/af/list.nhn?page=3
233 https://movie.naver.com/movie/point/af/list.nhn?page=4
234 ...
235 ...
236 <class 'pandas.core.frame.DataFrame'>
237 RangeIndex: 1000 entries, 0 to 999
238 Data columns (total 3 columns):
239 #   Column  Non-Null Count  Dtype
240 ---  ---
241 0   Title   1000 non-null    object
242 1   Point   1000 non-null    object
243 2   Review  1000 non-null    object
244 dtypes: object(3)
245 memory usage: 23.6+ KB
246
247
248
249 6. Coupang의 상품정보 Scraping
250 1)Web 문서들은 서로 다양한 문서 구조로 출력된다.
251 2)따라서 Python에서 web 문서로부터 scraping을 하기 위해서는 추출하고자 하는 정보들이 구성되어 있는 영역을 먼저 확인해야 한다.
252 3)Social Commerce의 대표적인 online market인 Coupang의 상품 정보 추출을 해보자.

```

253 4)`여성패션' 중 `여성 크로스백' 목록 item을 살펴보자.
254 5)Scraping하려는 web page의 URL 구조와 문서 구조를 파악한다.
255 6)URL 구조
256 <http://www.coupang.com/np/search?q=여성크로스백>
257 7)문서 구조
258 -상품명 : class="title"
259 -가격 : class="discount-price"
260
261
262
263 7. 한국일보 headline 기사 Scraping하기
264 1)한국일보 첫 page의 기사를 Scraping 해보자.
265 2)먼저 scraping 하려는 web page의 URL 구조와 문서 구조를 파악해야 한다.
266 3)URL 구조
267 <http://www.hankookilbo.com/>
268 4)문서 구조
269 -기사 제목 : class="title"