

```

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 2. 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=l!='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');</script>
100 <!-- End Google Tag Manager -->
101 ...
102 ...
103

```

```

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

```

```

107 리얼월드 HTTP : 역사와 코드로 배우는 인터넷과 웹 기술
108 이것이 MariaDB다
109 제프리 리처의 Windows via C/C++(북간판)
110 초보자를 위한 유니티 입문(개정판) : 따라 하면서 배우는 2D & 3D 게임 개발
111 회사에서 바로 통하는 실무 엑셀
112 맛있는 디자인 포토샵 CC 2019
113 알고리즘이 욕망하는 것들
114 파이썬 라이브러리를 활용한 머신러닝(번역개정판) : 사이킷런 핵심 개발자가 쓴 머신러닝과 데이터 과학 실무서
115 파이썬으로 웹 크롤러 만들기(2판) : 초간단 나만의 웹 크롤러로 원하는 데이터를 가져오는 방법
116 더 나은 세상을 위한 소프트 디지털
117 비도클래스 하원의 유튜브 동영상 편집 with 프리미어 프로
118 회사에서 바로 통하는 실무 엑셀+파워포인트+워드&한글
119 맛있는 디자인 포토샵&일러스트레이터 CC 2019
120 맛있는 디자인 프리미어 프로&애프터 이펙트 CC 2019
121 밑바닥부터 시작하는 딥러닝
122 이것이 C#이다
123 핸즈온 머신러닝
124 소문난 명강의 : 레트로의 유니티 게임 프로그래밍 에센스
125 이것이 자바다
126 이것이 우분투 리눅스다
127
128
129

```

3. Naver 영화 평점 Scraping 하기

```

130 from bs4 import BeautifulSoup
131
132 html_data = requests.get('https://movie.naver.com/movie/point/af/list.nhn?page=1')
133 soup = BeautifulSoup(html_data.text, 'html.parser')
134 titles = soup.find_all(class_='movie')
135
136 title_list = []
137 for title in titles:
138     print(title.text)
139 -----
140 대학살의 신
141 스타워즈: 라스트 제다이
142 내안의 그놈
143 말모이
144 말모이
145 내안의 그놈
146 연니
147 내안의 그놈
148 존 워
149 마이 리틀 자이언트
150
151
152 for title in titles:
153     title_list.append(title.text)
154
155 point_list = []
156 points = soup.find_all(class_='point')
157 for point in points:
158     point_list.append(point.text)
159
160 review_list = []
161 reviews = soup.find_all(class_='title')
162 for review in reviews:
163     rev = review.text
164     rev = rev.strip()
165     rev = rev.replace('\t', '')
166     rev = rev.replace('\n', '')
167     rev = rev.replace('신고', '')
168

```

```

169         review_list.append(rev)
170
171 df = pd.DataFrame(title_list, columns=['Title'])
172 df['Point'] = point_list
173 df['Review'] = review_list
174 df
175 -----
176      Title      Point  Review
177 0 대학살의 신      7      대학살의 신자식싸움에 부모등 터진다
178 1 스타워즈: 라스트 제다이      1      스타워즈: 라스트 제다이어거 보느니 로그원 열번 보는게 낫다
179 2 내안의 그놈      10      내안의 그놈재밌어요 유치해도 뽕뽕터짐 진영 연기 잘하네요 라미란과 잘...
180 3 말모이      10      말모이후반부에 보고 올었습니다 감동적임
181 4 말모이      10      말모이재미를 떠나서 역사는 무조건 10점이다
182 5 내안의 그놈      10      내안의 그놈뽕한스토리이지만 재밌게 잘 보고왔어요최고의성형은 다이어트
183 6 언니      5      언니론다 로우지가 언니 역활했으면 그나마 공감이 됐을듯...개연성도 떨어지고 액션도...
184 7 내안의 그놈      9      내안의 그놈ㅋㅋ재밌었음 생각보다 안정적인 연기력 소소하게 웃기 좋은 영화
185 8 존 워      10      존 워액션의 선두주자 키아누리브스!!
186 9 마이 리틀 자이언트      7      마이 리틀 자이언트동화보다 더 환상적인 거인
187
188
189

```

4. Naver 평점 1page부터 100page까지 scraping 하기

```

191 from bs4 import BeautifulSoup
192 import requests
193 import pandas as pd
194
195 url = 'https://movie.naver.com/movie/point/af/list.nhn?page='
196
197
198 title_list = []
199 point_list = []
200 review_list = []
201
202 for pge in range(1, 101):
203     url = url + str(pge)
204     print(url)
205     html_data = requests.get(url)
206     soup = BeautifulSoup(html_data.text, 'html.parser')
207     titles = soup.find_all(class_='movie')
208     points = soup.find_all(class_='point')
209     reviews = soup.find_all(class_='title')
210     for title in titles:
211         title_list.append(title.text)
212     for point in points:
213         point_list.append(point.text)
214     for review in reviews:
215         rev = review.text
216         rev = rev.strip()
217         rev = rev.replace('\t', '')
218         rev = rev.replace('\n', '')
219         rev = rev.replace('신고', '')
220         review_list.append(rev)
221     url = url.split('=')[0] + '='
222
223 df = pd.DataFrame(title_list, columns=['Title'])
224 df['Point'] = point_list
225 df['Review'] = review_list
226

```

```
df.info()
```

```

227 -----
228
229 <class 'pandas.core.frame.DataFrame'>
230 RangeIndex: 1020 entries, 0 to 1019
231 Data columns (total 3 columns):
232 Title    1020 non-null object
233 Point    1020 non-null object
234 Review   1020 non-null object
235 dtypes: object(3)
236 memory usage: 24.0+ KB
237
238
239

```

5. Coupang의 상품정보 Scraping

```

241 1)Web 문서들은 서로 다양한 문서 구조로 출력된다.
242 2)따라서 Python에서 web 문서로부터 scraping을 하기 위해서는 추출하고자 하는 정보들이 구성되어 있는 영역을 먼저 확인해야 한다.
243 3)Social Commerce의 대표적인 online market인 Coupang의 상품 정보 추출을 해보자.
244 4)'여성패션' 중 '여성 크로스백' 목록 item을 살펴보자.
245 5)Scraping하려는 web page의 URL 구조와 문서 구조를 파악한다.
246 6)URL 구조
247     http://www.coupang.com/np/search?q=여성크로스백
248 7)문서 구조
249     -상품명 : class="name"
250     -가격 : class="price-value"
251
252

```

253
254 6. 한국일보 headline 기사 Scraping하기
255 1)한국일보 첫 page의 기사를 Scraping 해보자.
256 2)먼저 scraping 하려는 web page의 URL 구조와 문서 구조를 파악해야 한다.
257 3)URL 구조
258 http://www.hankookilbo.com/
259 4)문서 구조
260 -기사 제목 : class="title"