

```

1 <준비>
2 1. Menu 중 Terminal > New Terminal > Python Version 확인
3 2. Ctrl + Shift + P --> 명령창에서 exe 엔터
4 3. 확장에서 Python, Python for VSCode, Python Extension Pack, Python(PyDev) 설치할 것
5 4. Google에서 'selenium python'으로 검색
6 5. https://www.seleniumhq.org/download/에서 Browser 중 Google Chrome Driver 설치할 것
7
8 1. Web Scraping이란?
9 1)Web 문서로부터 필요한 정보만을 추출하여 제공하는 기술을 말한다.
10 2)이 기술을 통해 상품 Catalog를 제작하거나 News 기사, Blog나 Cafe의 게시물, 회사의 profile과 금융 data 등을 수집할 수 있다.
11 3)그러기 위해서는 먼저 추출하고자 하는 정보들이 구성되어 있는 영역을 확인해야 한다.
12
13 2. Internet에서 data 수집하기
14 1)Data 수집은 data 준비 절차의 첫 단계로 어떤 형식의 data를 수집할지, 어떤 방법과 경로로 data를 수집할 것인가를 고민해야 한다.
15 2)수집된 data의 품질은 data 분석 결과에 많은 영향을 미치며, data가 얼마나 잘 정제되어 있는지에 따라서도 전체 data 분석에 드는 시간과 노력이 좌우되기 때문에 data 수집 과정은 매우 중요한 단계라 할 수 있다.
16
17 3. Web site에서 사용된 web 기술 확인
18 1)해당 web site에 사용된 web 기술을 확인하는 유용한 도구는 builtwith module이다.
19 2)Install
20 $ pip install builtwith
21 3)이 module은 전달된 URL을 가진 web site를 download하고 분석할 것이다.
22 4)분석이 되면 web site에 사용된 기술을 알려준다.
23
24 import builtwith
25 builtwith.parse('http://example.webscraping.com')
26 -----
27 {'web-servers': ['Nginx'],
28  'web-frameworks': ['Web2py', 'Twitter Bootstrap'],
29  'programming-languages': ['Python'],
30  'javascript-frameworks': ['jQuery', 'Modernizr', 'jQuery UI']}
31
32 4. Web site 소유자 찾기
33 1)website의 소유자를 찾으려면 해당 website의 domain명에 누가 등록됐는지 확인하는 whois protocol을 사용할 수 있다.
34 2)Install
35 $ pip install python-whois
36
37 import whois
38 print(whois.whois('appspot.com'))
39 -----
40 {'domain_name': ['APPSPOT.COM', 'appspot.com'],
41  'registrar': 'MarkMonitor, Inc.',
42  'whois_server': 'whois.markmonitor.com',
43  'referral_url': None,
44  'updated_date': [datetime.datetime(2019, 2, 6, 10, 33, 49),
45                  datetime.datetime(2019, 2, 6, 2, 33, 49)],
46  'creation_date': [datetime.datetime(2005, 3, 10, 2, 27, 55),

```

```

47     datetime.datetime(2005, 3, 9, 18, 27, 55)],
48     'expiration_date': [datetime.datetime(2020, 3, 10, 1, 27, 55),
49     datetime.datetime(2020, 3, 9, 0, 0)],
50     'name_servers': ['NS1.GOOGLE.COM',
51     'NS2.GOOGLE.COM',
52     'NS3.GOOGLE.COM',
53     'NS4.GOOGLE.COM',
54     'ns1.google.com',
55     'ns4.google.com',
56     'ns3.google.com',
57     'ns2.google.com'],
58     'status': ['clientDeleteProhibited https://icann.org/epp#clientDeleteProhibited',
59     'clientTransferProhibited https://icann.org/epp#clientTransferProhibited',
60     'clientUpdateProhibited https://icann.org/epp#clientUpdateProhibited',
61     'serverDeleteProhibited https://icann.org/epp#serverDeleteProhibited',
62     'serverTransferProhibited https://icann.org/epp#serverTransferProhibited',
63     'serverUpdateProhibited https://icann.org/epp#serverUpdateProhibited',
64     'clientUpdateProhibited (https://www.icann.org/epp#clientUpdateProhibited)',
65     'clientTransferProhibited (https://www.icann.org/epp#clientTransferProhibited)',
66     'clientDeleteProhibited (https://www.icann.org/epp#clientDeleteProhibited)',
67     'serverUpdateProhibited (https://www.icann.org/epp#serverUpdateProhibited)',
68     'serverTransferProhibited
69     (https://www.icann.org/epp#serverTransferProhibited)',
70     'serverDeleteProhibited (https://www.icann.org/epp#serverDeleteProhibited)'],
71     'emails': ['abusecomplaints@markmonitor.com',
72     'whoisrequest@markmonitor.com'],
73     'dnssec': 'unsigned',
74     'name': None,
75     'org': 'Google LLC',
76     'address': None,
77     'city': None,
78     'state': 'CA',
79     'zipcode': None,
80     'country': 'US'}

```

5. Website crawling

- 1) Web site를 scrap을 하려면 먼저 crawling 이라고 알려진 과정, 즉 관심 있는 data를 가진 web page를 download할 필요가 있다.
- 2) Web site crawling을 할 수 있는 몇 가지 방법이 있으며, 대상 web site 구조에 맞춰 적절히 선택한다.
- 3) Web site download
 - Web page를 crawling하려면 우선 web page를 download 해야 한다.
 - 다음은 web 주소가 URL인 web site를 download하고자 Python의 urllib3 module을 사용하는 간단한 Python script이다.

```

86     from urllib.request import urlopen
87
88     def download(url):
89         return urlopen(url).read().decode('utf-8')

```

- 이 함수는 URL이 전달됐을 때 web page를 download하고 HTML을 반환한다.
- 이 source code는 web page를 download할 때 처리가 안되는 오류에 직면할 수도 있는 문제

가 있다.

- 이를 테면 요청한 page가 존재하지 않는 것이 그런 사례이다.
- 그럴 경우 urllib3는 예외를 발생한 후 script의 실행을 멈춘다.
- 그럴 경우를 대비해서 아래의 code로 변경하자.

```

98 from urllib.request import urlopen
99 from urllib.error import HTTPError
100
101 def download(url):
102     print('downloading:', url)
103     try :
104         html = urlopen(url).read().decode('utf-8')
105     except HTTPError as e:
106         print('Download error:', e.reason)
107         html = None
108     return html

```

- 이제 download 오류 상황이 되면 예외 처리가 되고 None을 반환한다.
- 이 기능을 test하려면 500 오류를 반환하는 <http://httpstat.us/500> page를 download하면 된다.

```

112 download('http://httpstat.us/500')
113 -----
114 Downloading: http://httpstat.us/500
115 Download error: Internal Server Error
116
117 download('http://www.samsung.com')
118 -----
119 Downloading: http://www.samsung.com
120 Download error: Forbidden

```

4)Download 재시도

- Download할 때 가끔 발생하는 오류가 일시적인 것들도 있다.
- 예를 들어, Web server가 과부하되면 503 Service Unavailable 오류가 발생한다.
- 이 오류는 server 상의 문제가 해결되면 더 이상 나타나지 않고 web page를 다시 download할 수 있다.
- 그러나 모든 오류에 대해 download를 다시 시도하지 않으려 한다.
- Server에서 404 Not Found 오류가 발생하면 web page가 현재 존재하지 않기 때문에 계속 같은 오류가 발생할 것이다.
- 일어날 가능성이 있는 HTTP 오류 전체 list는 IETF(Internet Engineering Task Force)에서 정의했고, <http://tools.ietf.org/html/rfc7231#section-6>에서 확인할 수 있다.
- 이 문서에서 4xx 오류들은 요청할 때 잘못된 것이 있을 때 일어나며, 5xx 오류들은 server에 문제가 있을 때 일어난다는 것을 알 수 있다.
- 따라서 5xx 오류가 일어날 때만 download를 재시도하게 download 함수를 만들자.

```

134
135 from urllib.request import urlopen
136 from urllib.error import HTTPError
137
138 def download(url, num_retries=2):
139     print('downloading:', url)

```

```

140     try :
141         html = urlopen(url).read().decode('utf-8')
142     except HTTPError as e:
143         print('Download error:', e.reason)
144         html = None
145         if num_retries > 0:
146             if hasattr(e, 'code') and 500 <= e.code < 600:
147                 #5xx HTTP 오류시 재 시도
148                 return download(url, num_retries - 1)
149     return html
150
151 -이제 5xx 오류가 발생하면 download 함수는 계속 자신을 호출하면서 재시도한다.
152 -또한 download 함수는 재시도할 수 있는 횟수를 알려주는 전달 인자도 갖고 있다.
153 -이 전달 인자는 기본 2회로 되어 있다.
154 -Server 오류가 해결되지 않을 수도 있기 때문에 web page를 download하는 시도 횟수를 제한
155 한다.
156 -이 기능을 test하려면 500 오류를 반환하는 http://httpstat.us/500 page를 download하면
157 된다.
158
159 download('http://httpstat.us/500')
160 downloading: http://httpstat.us/500
161 Download error: Internal Server Error
162 downloading: http://httpstat.us/500
163 Download error: Internal Server Error
164 downloading: http://httpstat.us/500
165 Download error: Internal Server Error
166
167 -예상대로 download 함수는 web page를 download하려 한다.
168 -그리고 오류가 발생하면 함수를 끝내기 전에 2번 더 download를 시도한다.
169
170 6. Crawling & Scraping
171 1)Crawling
172 -Web Site의 data를 그대로 취득하는 것
173 2)Scraping
174 -Crawling하여 모든 data에서 필요한 것만 추출하거나 변환하는 처리를 포함.
175 3)정규식
176 -정규식은 변경 사항에 대해 유연하게 대처할 수는 있지만 만들기 어렵고 가독성도 떨어진다.
177 -따라서 적용하기에 너무 취약하고 web page가 바뀌면 사용하는 정규식이 쉽게 무용지물이 된다는
178 점의 문제점이 있다.
179
180 4)pandas.read_html() 함수 사용하기
181 -다음의 third-party package들의 설치여부를 확인한다.
182 --html5lib, lxml, BeautifulSoup4
183 -HTML file의 table 요소를 DataFrame으로 불러온다.
184 -DataFrame이 들어 있던 list가 return된다.
185 -Table 요소가 여러 개 있는 경우, 여러 개의 DataFrame이 저장된다.
186
187 import pandas as pd
188
189 url = 'http://www.fdic.gov/bank/individual/failed/banklist.html'
190
191 dfs = pd.read_html(url)

```

```

189
190     dfs[0].info()
191     -----
192     -
193     <class 'pandas.core.frame.DataFrame'>
194     RangeIndex: 555 entries, 0 to 554
195     Data columns (total 7 columns):
196     Bank Name          555 non-null object
197     City                555 non-null object
198     ST                  555 non-null object
199     CERT                555 non-null int64
200     Acquiring Institution 555 non-null object
201     Closing Date        555 non-null object
202     Updated Date        555 non-null object
203     dtypes: int64(1), object(6)
204     memory usage: 30.4+ KB
205
206 5)requests module 이용하기
207   -Python HTTP for Humans'
208   -http://docs.python-requests.org/en/master/
209   -Install
210     pip install requests
211
212     import requests
213     url = 'https://www.naver.com'
214     naver = requests.get(url)
215     print(naver.text)
216
217   -requests module을 사용하여 Web API로 Wikipedia의 정보를 얻는 방법을 소개한다.
218   -이 Wikipedia의 data를 program으로 획득하기 위해 MediaWiki라는 service를 사용할 것이
219     다.
220     --https://www.mediawiki.org/wiki/MediaWiki
221
222     import requests, pprint
223     api_url = 'https://en.wikipedia.org/w/api.php'
224     api_params = {'format':'json', 'action':'query', 'titles':'Jack
225     Bauer','prop':'revisions', 'rvprop':'content'}
226     #title : 검색하고 싶은 keyword
227     #prop : 검색 결과로 어떤 정보를 반환할지를 지정
228     #rvprop : prop으로 지정한 항목을 더 구체적으로 지정할 수있다.
229
230     wiki_data = requests.get(api_url, params=api_params).json()
231     pprint.pprint(wiki_data)
232     -----
233     {'batchcomplete': '',
234      'query': {'pages': {'389903': {'ns': 0,
235                                   'pageid': 389903,
236                                   'revisions': [{'*': '{Use mdy '
237                                   'dates|date=December '
238                                   '2017}}\n'
239                                   '{{other people}}\n'
240                                   '{{Infobox character\n'

```

```

238         '| color          = '
239         '#B3001B\n'
240         '| name           = '
241         'Jack Bauer\n'
242         '| series         = '
243         '[[24 (TV series)|24]]\n'
244         '| image          = '
245         'Jack Bauer.jpg\n'
246         '| caption        = '
247     ...
248     ...

```

-action에 지정할 수 있는 것

<https://en.wikipedia.org/w/api.php>

-action으로 query를 지정했을 때 prop에 지정할 수 있는 것

<https://en.wikipedia.org/w/api.php?action=help&modules=query>

-prop로 revisions를 지정했을 때 rvprop에 지정할 수 있는 것

<https://en.wikipedia.org/w/api.php?action=help&modules=query%2Brevisions>

-다양한 조건을 지정할 수 있지만, wikipedia 검색 결과를 표시하고 싶다면 위의 예제에서 지정한 설정을 그대로 사용하고, title parameter만 바꾸면 된다.

```

260
261     import requests
262     import codecs
263     api_base_url = 'https://en.wikipedia.org/w/api.php'
264     api_params = {'format':'xmlfm', 'action':'query', 'titles':'Jack
265                 Bauer','prop':'revisions', 'rvprop':'content'}
266     wiki_data = requests.get(api_base_url, params=api_params)
267     fo = codecs.open('C:/temp/wiki.html', 'w', 'utf-8')
268     fo.write(wiki_data.text)
269     fo.close()

```

6) pprint module 이용하기

-대량의 data를 보기 쉽게 표시해주는 표준 module

-pprint(prettyprint)

```

270
271
272
273
274     import requests
275     import pprint
276     url = 'https://www.naver.com'
277     naver = requests.get(url)
278     pprint.pprint(naver.text)
279

```

7) BeautifulSoup module 이용하기

-Scraping을 위한 module

-일반적으로 HTML tag들이 start tag와 end tag가 서로 pair 되지 않을 경우가 많다.

-pair 되지 않아도 아름답게 처리해 주는 module

<https://www.crummy.com/software/BeautifulSoup/bs4/doc/>

```

280
281
282
283
284
285
286     from bs4 import BeautifulSoup
287     import requests

```

```

288
289     html_data = requests.get('https://www.naver.com')
290     soup = BeautifulSoup(html_data.text, 'html.parser')
291     soup.title
292     -----
293     <title>NAVER</title>
294
295 [한빛미디어 책 제목 읽어오기]
296     hanbit = requests.get('http://www.hanbit.co.kr/media/')
297     soup = BeautifulSoup(hanbit.text, 'html.parser')
298     soup
299     -----
300     <!DOCTYPE html>
301     <html lang="ko">
302     <head>
303     <!--[if lte IE 8]>
304     <script>
305         location.replace('/support/explorer_upgrade.html');
306     </script>
307     <![endif]-->
308     <!-- Google Tag Manager -->
309     <script>(function(w,d,s,l,i){w[l]=w[l]||[];w[l].push({'gtm.start':
310     new Date().getTime(),event:'gtm.js'});var f=d.getElementsByTagName(s)[0],
311     j=d.createElement(s),dl=l!='dataLayer'?'&l='+l:'';j.async=true;j.src=
312     'https://www.googletagmanager.com/gtm.js?id='+i+dl;f.parentNode.insertBefore(Befo
313     re(j,f);
314     })(window,document,'script','dataLayer','GTM-W9D5PM3');</script>
315     <!-- End Google Tag Manager -->
316     ...
317     ...
318     for book in soup.find_all('p', class_='book_tit'):
319         print(book.find('a').text)
320     -----
321     리얼월드 HTTP : 역사와 코드로 배우는 인터넷과 웹 기술
322     이것이 MariaDB다
323     제프리 리처의 Windows via C/C++(북간판)
324     초보자를 위한 유니티 입문(개정판) : 따라 하면서 배우는 2D & 3D 게임 개발
325     회사에서 바로 통하는 실무 엑셀
326     맛있는 디자인 포토샵 CC 2019
327     알고리즘이 욕망하는 것들
328     파이썬 라이브러리를 활용한 머신러닝(번역개정판) : 사이킷런 핵심 개발자가 쓴 머신러닝과 데이
329     터 과학 실무서
330     파이썬으로 웹 크롤러 만들기(2판) : 초간단 나만의 웹 크롤러로 원하는 데이터를 가져오는 방법
331     더 나은 세상을 위한 소프트 디지털
332     비도클래스 하천의 유튜브 동영상 편집 with 프리미어 프로
333     회사에서 바로 통하는 실무 엑셀+파워포인트+워드&한글
334     맛있는 디자인 포토샵&일러스트레이터 CC 2019
335     맛있는 디자인 프리미어 프로&애프터 이펙트 CC 2019
336     밑바닥부터 시작하는 딥러닝
337     이것이 C#이다
338     핸즈온 머신러닝

```



```

338     소문난 명강의 : 레트로의 유니티 게임 프로그래밍 에센스
339     이것이 자바다
340     이것이 우분투 리눅스다
341
342 [Naver 영화 평점 Scraping 하기]
343
344     from bs4 import BeautifulSoup
345
346     html_data =
347     requests.get('https://movie.naver.com/movie/point/af/list.nhn?page=1')
348     soup = BeautifulSoup(html_data.text, 'html.parser')
349     titles = soup.find_all(class_='movie')
350
351     title_list = []
352     for title in titles:
353         print(title.text)
354         -----
355         대학살의 신
356         스타워즈: 라스트 제다이
357         내안의 그놈
358         말모이
359         말모이
360         내안의 그놈
361         언니
362         내안의 그놈
363         존 워
364         마이 리틀 자이언트
365
366     for title in titles:
367         title_list.append(title.text)
368
369     point_list = []
370     points = soup.find_all(class_='point')
371     for point in points:
372         point_list.append(point.text)
373
374     review_list = []
375     reviews = soup.find_all(class_='title')
376     for review in reviews:
377         rev = review.text
378         rev = rev.strip()
379         rev = rev.replace('\t', '')
380         rev = rev.replace('\n', '')
381         rev = rev.replace('신고', '')
382         review_list.append(rev)
383
384     df = pd.DataFrame(title_list, columns=['Title'])
385     df['Point'] = point_list
386     df['Review'] = review_list
387     df
388     -----
389     Title                Point    Review

```


389	0 대학살의 신	7	대학살의 신자식싸움에 부모등 터진다
390	1 스타워즈: 라스트 제다이	1	스타워즈: 라스트 제다이어거 보느니 로그원 열번 보는게 낫다
391	2 내안의 그놈	10	내안의 그놈재밌어요 유치해도 뽕뽕터짐 진영 연기 잘하네요 라미란과 잘...
392	3 말모이	10	말모이후반부에 보고 울었습니다 감동적임
393	4 말모이	10	말모이재미를 떠나서 역사는 무조건 10점이다
394	5 내안의 그놈	10	내안의 그놈뽕한스토리이지만 재밌게 잘 보고왔어요최고 의성형은 다이어트
395	6 언니	5	언니론다 로우지가 언니 역활했으면 그나마 공감이 됐을 듯...개연성도 떨어지고 액션도...
396	7 내안의 그놈	9	내안의 그놈ㅋㅋ재밌었음 생각보다 안정적인 연기력 소소하게 웃기 좋은 영화
397	8 존 워	10	존 워액션의 선두주자 키아누리브스!!
398	9 마이 리틀 자이언트	7	마이 리틀 자이언트동화보다 더 환상적인 거인

[Naver 평점 1page부터 100page까지 scraping 하기]

```

401
402 from bs4 import BeautifulSoup
403
404 url = 'https://movie.naver.com/movie/point/af/list.nhn?page='
405
406 title_list = []
407 point_list = []
408 review_list = []
409
410 for pge in range(1, 101):
411     url = url + str(pge)
412     print(url)
413     html_data = requests.get(url)
414     soup = BeautifulSoup(html_data.text, 'html.parser')
415     titles = soup.find_all(class_='movie')
416     points = soup.find_all(class_='point')
417     reviews = soup.find_all(class_='title')
418     for title in titles:
419         title_list.append(title.text)
420     for point in points:
421         point_list.append(point.text)
422     for review in reviews:
423         rev = review.text
424         rev = rev.strip()
425         rev = rev.replace('\t', '')
426         rev = rev.replace('\n', '')
427         rev = rev.replace('신고', '')
428         review_list.append(rev)
429     url = url.split('=')[0] + '='
430
431 df = pd.DataFrame(title_list, columns=['Title'])
432 df['Point'] = point_list
433 df['Review'] = review_list
434
435 df.info()

```

```

436 -----
437 <class 'pandas.core.frame.DataFrame'>
438 RangeIndex: 1020 entries, 0 to 1019
439 Data columns (total 3 columns):
440 Title    1020 non-null object
441 Point    1020 non-null object
442 Review   1020 non-null object
443 dtypes: object(3)
444 memory usage: 24.0+ KB
445
446 [Coupang의 상품정보 Scraping]
447 -Web 문서들은 서로 다양한 문서 구조로 출력된다.
448 -따라서 R에서 web 문서로부터 scraping을 하기 위해서는 추출하고자 하는 정보들이 구성되어
    있는 영역을 먼저 확인해야 한다.
449 -Social Commerce의 대표적인 online market인 Coupang의 상품 정보 추출을 해보자.
450 -'여성패션' 중 '여성 크로스백' 목록 item을 살펴보자.
451 -Scraping하려는 web page의 URL 구조와 문서 구조를 파악한다.
452 -URL 구조
453 http://www.coupang.com/np/search?q=여성크로스백
454 -문서 구조
455   --상품명 : class="name"
456   --가격 : class="price-value"
457
458 [한국일보 headline 기사 Scraping하기]
459 -한국일보 첫 page의 기사를 Scraping 해보자.
460 -먼저 scraping 하려는 web page의 URL 구조와 문서 구조를 파악해야 한다.
461 -URL 구조
462 http://www.hankookilbo.com/
463 -문서 구조
464   --기사 제목 : class="firstList"
465
466
467 8)lxml module 이용하기
468 -이 module은 C로 만들어진 libxml2 XML 분석 library를 Python에서 사용할 수 있게 만든
    module이다.
469 -Beautiful Soup보다 더 빠르게 분석하고 scraping할 수 있지만, 어떤 computer에서는 설치가
    좀 어려운 편이다.
470 -Install
471   $ pip install lxml
472 -필요하다면 cssselect 도 install
473   $ pip install cssselect
474
475   import lxml.html
476
477   url = 'http://www.hanbit.co.kr/media/'
478   html = download(url)
479   tree = lxml.html.fromstring(html)
480   for a in tree.cssselect('a'):
481       #href 속성과 글자를 추출한다.
482       print(a.get('href'), a.text)
483
484 9)requests를 사용하여 API에 접근하기

```

```

485 -기상청 RSS(http://www.weather.go.kr/weather/lifenindustry/sevice\_rss.jsp)를 이용
486 하자.
487 -2013년 동네 예보 RSS
488 http://www.weather.go.kr/images/weather/lifenindustry/dongnaeforecast\_rss.p
489 df
490 -주소선택후 rss button click하면 zone을 알 수 있다.
491 --예:서울특별시 강남구 청담동
492 http://www.kma.go.kr/wid/queryDFSRSS.jsp?zone=1168056500
493
494 import requests
495 import pprint
496 api_uri = 'http://www.kma.go.kr/wid/queryDFSRSS.jsp?zone=1168056500'
497 weather_data = requests.get(api_uri).text
498 pprint.pprint(weather_data)
499 -----
500 ('<?xml version="1.0" encoding="UTF-8" ?>\n'
501  '<rss version="2.0">\n'
502  '<channel>\n'
503  '<title>기상청 동네예보 웹서비스 - 서울특별시 강남구 청담동 도표예보</title>\n'
504  '<link>http://www.kma.go.kr/weather/main.jsp</link>\n'
505  '<description>동네예보 웹서비스</description>\n'
506  '<language>ko</language>\n'
507  '<generator>동네예보</generator>\n'
508  '<pubDate>2019년 01월 15일 (화)요일 14:00</pubDate>\n'
509  '<item>\n'
510  '<author>기상청</author>\n'
511  '<category>서울특별시 강남구 청담동</category>\n'
512  '<title>동네예보(도표) : 서울특별시 강남구 청담동 '
513  '[X=61,Y=126]</title><link>http://www.kma.go.kr/weather/forecast/timeseries.jsp?searchType=INTEREST&dongCode=1168056500</link>\n'
514  '<guid>http://www.kma.go.kr/weather/forecast/timeseries.jsp?searchType=INTEREST&dongCode=1168056500</guid>\n'
515  '<description>\n'
516  ...
517  ...
518 -get method의 params option을 활용하기
519
520 api_url = 'http://www.kma.go.kr/wid/queryDFSRSS.jsp'
521
522 payload = {'zone':'1168056500'}
523
524 weather_data = requests.get(api_url, payload).text
525
526 weather_data
527 -----
528 '<?xml version="1.0" encoding="UTF-8" ?>\n<rss
529 version="2.0">\n<channel>\n<title>기상청 동네예보 웹서비스 - 서울특별시 강남구 청
530 담동 도표예보
531 </title>\n<link>http://www.kma.go.kr/weather/main.jsp</link>\n<description
532 >동네예보 웹서비스</description>\n<language>ko</language>\n<generator>동네
533 예보</generator>\n<pubDate>2019년 01월 15일 (화)요일 14:00</pubDate>\n

```

```
<item>\n<author>기상청</author>\n<category>서울특별시 강남구 청담동
</category>\n<title>동네예보(도표) : 서울특별시 강남구 청담동
[X=61,Y=126]</title><link>http://www.kma.go.kr/weather/forecast/timeserie
s.jsp?searchType=INTEREST&dongCode=1168056500</link>\n<guid>http://www.kma.go.kr/weather/forecast/timeseries.jsp?searchType=INTEREST&a
mp;dongCode=1168056500</guid>\n<description>\n <header>\n
<tm>201901151400</tm>\n <ts>4</ts>\n
```

10)xml.etree.ElementTree module 사용하기
-Python에서 xml data를 다루기 위한 module

```
import xml.etree.ElementTree as ET
```

```
xml_data = ET.fromstring(weather_data)
```

```
for tag in xml_data.iter('data'):
    print(tag.find('hour').text + "/" + tag.find('temp').text)
```

```
-----
18/-2.0
```

```
21/-4.0
```

```
24/-5.0
```

```
...
```

```
...
```

```
list = []
```

```
for tag in xml_data.iter('data'):
    dic = {'hour': tag.find('hour').text,
          'day': tag.find('day').text,
          'temp': tag.find('temp').text,
          'tmx': tag.find('tmx').text,
          'tmn': tag.find('tmn').text,
          'sky': tag.find('sky').text,
          'pty': tag.find('pty').text,
          'wfKor': tag.find('wfKor').text,
          'wfEn': tag.find('wfEn').text}
```

```
list.append(dic)
```

```
df = pd.DataFrame(list, columns=['hour', 'day', 'temp', 'tmx', 'tmn', 'sky', 'pty',
'wfKor', 'wfEn'])
```

```
df
```

```
-----
   hour  day temp   tmx   tmn  sky  pty wfKor  wfEn
565  0 18  0  -2.0 -999.0 -999.0  2   0   구름 조금 Partly Cloudy
566  1 21  0  -4.0 -999.0 -999.0  2   0   구름 조금 Partly Cloudy
567  2 24  0  -5.0 -999.0 -999.0  1   0   맑음      Clear
568  3  3  1  -6.0  -1.0  -8.0  1   0   맑음      Clear
569  4  6  1  -7.0  -1.0  -8.0  1   0   맑음      Clear
```

11)RSS Scraping

```

572 -전자신문 RSS
573 import pandas as pd
574 import xml.etree.ElementTree as ET
575 import requests
576
577 api_url = 'http://rss.etnews.com/Section901.xml'
578
579 etnews_data = requests.get(api_url).text
580 etnews_data
581 -----
582 '<?xml version="1.0" encoding="utf-8" ?>\n<rss version="2.0">\r\n
<channel>\r\n  ...
583 ...
584
585 xml_data = ET.fromstring(etnews_data)
586 for tag in xml_data.iter('item'):
587     print(tag.find('title').text + ", " + tag.find('pubDate').text)
588 -----
589 애플 납품 업체들, 줄줄이 실적 하향... '아이폰 쇼크' 후폭풍,Tue, 15 Jan 2019 17:00:00
+0900
590 지난해 드론 자격증 취득자 전년비 4배 증가...실효성 지적도,Tue, 15 Jan 2019 17:00:00
+0900
591 화웨이, 韓 스마트워치 시장 진출,Tue, 15 Jan 2019 17:00:00 +0900
592 ...
593 ...
594
595 list = []
596
597 for tag in xml_data.iter('item'):
598     dic = {'Title':tag.find('title').text,
599           'Link':tag.find('link').text,
600           'Author':tag.find('author').text,
601           'PubDate':tag.find('pubDate').text,
602           'Guid':tag.find('guid').text}
603     list.append(dic)
604
605 list
606 -----
607 [{'Title': "애플 납품 업체들, 줄줄이 실적 하향... '아이폰 쇼크' 후폭풍",
608   'Link': 'http://www.etnews.com/20190115000273',
609   'Author': '윤건일',
610   'PubDate': 'Tue, 15 Jan 2019 17:00:00 +0900',
611   'Guid': '20190115000273'},
612  ...
613  ...
614
615 df = pd.DataFrame(list, columns=['Title', 'Link', 'Author', 'PubDate', 'Guid'])
616 df.head()
617 -----
618 ...
619 ...
620

```

```

621     df.info()
622     -----
623     <class 'pandas.core.frame.DataFrame'>
624     RangeIndex: 30 entries, 0 to 29
625     Data columns (total 5 columns):
626     Title      30 non-null object
627     Link       30 non-null object
628     Author     30 non-null object
629     PubDate    30 non-null object
630     Guid       30 non-null object
631     dtypes: object(5)
632     memory usage: 1.2+ KB
633
634 7. Web site의 data file 읽기
635 -Web site에 있는 data set를 R로 loading 하는 방법을 알아보자.
636 -다음은 Titanic data file을 CSV file로 저장하는 방법을 소개한다.
637 -아래의 site에 접속해 보자.
638 -https://github.com/vincentarelbundock/Rdatasets/tree/master/csv/datasets
639 -https://vincentarelbundock.github.io/Rdatasets/datasets.html
640 --Item(datasets에서 찾는다) 중 'Titanic'을 찾아보자.
641 --Link의 속성을 파악하기 위해서는 Google Chrome보다는 Internet Explorer를 이용하는 것
    이 좋다.
642 --Internet Explorer로 해당 site를 방문한 다음, datasets의 Item 중 CSV의 link를 Mouse
    오른 Click를 한 후, [속성]을 선택한다.
643 --속성 창에서 Data set file의 URL을 확인하자.
644
645 import pandas as pd
646 url = 'https://vincentarelbundock.github.io/Rdatasets/csv/datasets/Titanic.csv'
647 df = pd.read_csv(url)
648 df.head()
649
650 8. Selenium
651 1)run.py
652 # 인터파크 투어 사이트에서 여행지를 입력후 검색 -> 잠시후 -> 결과
653 # 로그인시 PC 웹 사이트에서 처리가 어려울 경우 -> 모바일 로그인 진입
654 # 모듈 가져오기
655 # pip install selenium
656 # pip install bs4
657 # pip install pymysql
658 from selenium import webdriver as wd
659 from bs4 import BeautifulSoup as bs
660 from selenium.webdriver.common.by import By
661 # 명시적 대기를 위해
662 from selenium.webdriver.support.ui import WebDriverWait
663 from selenium.webdriver.support import expected_conditions as EC
664 from DbMgr import DBHelper as Db
665 import time
666 from Tour import TourInfo
667
668 # 사전에 필요한 정보를 로드 => 디비혹스 쉘, 배치 파일에서 인자로 받아서 세팅
669 db = Db()
670 main_url = 'http://tour.interpark.com/'

```

```

671 keyword = '로마'
672 # 상품 정보를 담은 리스트 (TourInfo 리스트)
673 tour_list = []
674
675 # 드라이버 로드
676 # 맥용
677 # driver = wd.Chrome(executable_path='./chromedriver')
678 # 윈도우용
679 driver = wd.Chrome(executable_path='chromedriver.exe')
680 # 고스트용
681 # driver = wd.PhantomJS(executable_path='./phantomjs')
682 # 차후 -> 옵션 부여하여 (프록시, 에이전트 조작, 이미지를 배제)
683 # 크롤링을 오래돌리면 => 임시파일들이 쌓인다!! -> 템프 파일 삭제
684
685 # 사이트 접속 (get)
686 driver.get(main_url)
687 # 검색창을 찾아서 검색어 입력
688 # id : SearchGNBText
689 driver.find_element_by_id('SearchGNBText').send_keys(keyword)
690 # 수정할경우 => 뒤에 내용이 붙어버림 => .clear() -> send_keys('내용')
691 # 검색 버튼 클릭
692 driver.find_element_by_css_selector('button.search-btn').click()
693
694 # 잠시 대기 => 페이지가 로드되고 나서 즉각적으로 데이터를 획득 하는 행위는
695 # 명시적 대기 => 특정 요소가 로케이트(발결된때까지) 대기
696 try:
697     element = WebDriverWait(driver, 10).until(
698         # 지정한 한개 요소가 올라면 웨이트 종료
699         EC.presence_of_element_located( (By.CLASS_NAME, 'oTravelBox') )
700     )
701 except Exception as e:
702     print( '오류 발생', e)
703 # 암묵적 대기 => DOM이 다 로드 될때까지 대기 하고 먼저 로드되면 바로 진행
704 # 요소를 찾을 특정 시간 동안 DOM 풀링을 지시 예를 들어 10 초이내 라로
705 # 발견 되면 진행
706 driver.implicitly_wait( 10 )
707 # 절대기 대기 => time.sleep(10) -> 클라우드 페어(디도스 방어 솔루션)
708 # 더보기 눌러서 => 게시판 진입
709 driver.find_element_by_css_selector('.oTravelBox>.boxList>.moreBtnWrap>.moreBtn').click()
710
711 # 게시판에서 데이터를 가져올때
712 # 데이터가 많으면 세션(혹시 로그인을 해서 접근되는 사이트일 경우) 관리
713 # 특정 단위별로 로그아웃 로그인 계속 시도
714 # 특정 게시물이 사라질 경우 => 팝업 발생 (없는 ...) => 팝업 처리 검토
715 # 게시판 스캔시 => 임계점을 모름!!
716 # 게시판 스캔 => 메타 정보 획득 => loop 를 돌려서 일괄적으로 방문 접근 처리
717
718 # searchModule.SetCategoryList(1, "") 스크립트 실행
719 # 16은 임시값, 게시물을 넘어갔을때 현상을 확인차
720 for page in range(1, 2):#16):
721     try:

```



```

722     # 자바스크립트 구동하기
723     driver.execute_script("searchModule.SetCategoryList(%s, '')" % page)
724     time.sleep(2)
725     print("%s 페이지 이동" % page)
726     #####
727     # 여러 사이트에서 정보를 수집할 경우 공통 정보 정의 단계 필요
728     # 상품명, 코멘트, 기간1, 기간2, 가격, 평점, 썸네일, 링크(상품상세정보)
729     boxItems = driver.find_elements_by_css_selector('.oTravelBox>.boxList>li')
730     # 상품 하나 하나 접근
731     for li in boxItems:
732         # 이미지를 링크값을 사용할것인가?
733         # 직접 다운로드 해서 우리 서버에 업로드(ftp) 할것인가?
734         print( '썸네임', li.find_element_by_css_selector('img').get_attribute('src') )
735         print( '링크', li.find_element_by_css_selector('a').get_attribute('onclick') )
736         print( '상품명', li.find_element_by_css_selector('h5.proTit').text )
737         print( '코멘트', li.find_element_by_css_selector('.proSub').text )
738         print( '가격', li.find_element_by_css_selector('.proPrice').text )
739         area = ""
740         for info in li.find_elements_by_css_selector('.info-row .proInfo'):
741             print( info.text )
742         print('='*100)
743         # 데이터 모음
744         # li.find_elements_by_css_selector('.info-row .proInfo')[1].text
745         # 데이터가 부족하거나 없을수도 있으므로 직접 인덱스로 표현은 위험성이 있음
746         obj = TourInfo(
747             li.find_element_by_css_selector('h5.proTit').text,
748             li.find_element_by_css_selector('.proPrice').text,
749             li.find_elements_by_css_selector('.info-row .proInfo')[1].text,
750             li.find_element_by_css_selector('a').get_attribute('onclick'),
751             li.find_element_by_css_selector('img').get_attribute('src')
752         )
753         tour_list.append( obj )
754     except Exception as e1:
755         print( '오류', e1 )
756
757     print( tour_list, len(tour_list) )
758     # 수집한 정보 개수를 루프 => 페이지 방문 => 콘텐츠 획득(상품상세정보) => 디비
759     for tour in tour_list:
760         # tour => TourInfo
761         print( type(tour) )
762         # 링크 데이터에서 실데이터 획득
763         # 분해
764         arr = tour.link.split(',')
765         if arr:
766             # 대체
767             link = arr[0].replace('searchModule.OnClickDetail(',')')
768             # 슬라이싱 => 앞에 ', 뒤에 ' 제거
769             detail_url = link[1:-1]
770             # 상세 페이지 이동 : URL 값이 완성된 형태인지 확인 (http~)
771             driver.get( detail_url )
772             time.sleep(2)

```

```

773     # pip install bs4
774     # 현재 페이지를 BeautifulSoup 의 DOM으로 구성
775     soup = bs( driver.page_source, 'html.parser')
776     # 현재 상세 정보 페이지에서 스케줄 정보 획득
777     data = soup.select('.tip-cover')
778     #print( type(data), len(data), type(data[0].contents) )
779     # 디비 입력 => pip install pymysql
780     # 데이터 sum
781     content_final = ""
782     for c in data[0].contents:
783         content_final += str(c)
784
785     # html 콘텐츠 데이터 전처리 (디비에 입력 가능토록)
786     import re
787     content_final = re.sub("","", content_final)
788     content_final = re.sub(re.compile(r'\r\n|\r|\n|\n\r+'), "\n", content_final)
789
790     print( content_final )
791     # 콘텐츠 내용에 따라 전처리 => data[0].contents
792     db.db_insertCrawlingData(
793         tour.title,
794         tour.price[:-1],
795         tour.area.replace('출발 가능 기간 : ',''),
796         content_final,
797         keyword
798     )
799
800     # 종료
801     driver.close()
802     driver.quit()
803     import sys
804     sys.exit()
805
806 2)Tour.py
807     # 상품 정보를 담는 클래스
808     class TourInfo:
809         # 멤버변수 (실제 컬럼보다는 작게 세팅했음)
810         title = ""
811         price = ""
812         area = ""
813         link = ""
814         img = ""
815         contents = ""
816         # 생성자
817         def __init__(self, title, price, area, link, img, contents=None ):
818             self.title = title
819             self.price = price
820             self.area = area
821             self.link = link
822             self.img = img
823             self.contents = contents
824

```

```
825 3)DbMgr.py
826 # 디비 처리, 연결, 해제, 검색어 가져오기, 데이터 삽입
827 import pymysql as my
828
829 class DBHelper:
830     """
831     멤버변수 : 커넥션
832     """
833     conn = None
834     """
835     생성자
836     """
837     def __init__(self):
838         self.db_init()
839     """
840     멤버 함수
841     """
842     def db_init(self):
843         self.conn = my.connect(
844             host='localhost',
845             user='root',
846             password='1234',
847             db='pythonDB',
848             charset='utf8',
849             cursorclass=my.cursors.DictCursor )
850
851     def db_free(self):
852         if self.conn:
853             self.conn.close()
854
855     # 검색 키워드 가져오기 => 웹에서 검색
856     def db_selectKeyword(self):
857         # 커서 오픈
858         # with => 닫기를 처리를 자동으로 처리해준다 => I/O 많이 사용
859         rows = None
860         with self.conn.cursor() as cursor:
861             sql = "select * from tbl_keyword;"
862             cursor.execute(sql)
863             rows = cursor.fetchall()
864             print(rows)
865         return rows
866
867     def db_insertCrawlingData(self, title, price, area, contents, keyword ):
868         with self.conn.cursor() as cursor:
869             sql = '''
870             insert into `tbl_crawlingdata`
871             (title, price, area, contents, keyword)
872             values( %s,%s,%s,%s,%s )
873             '''
874             cursor.execute(sql, (title, price, area, contents, keyword) )
875             self.conn.commit()
876
```

```
877     # 단독으로 수행시에만 작동 => 테스트코드를 삽입해서 사용
878     if __name__ == '__main__':
879         db = DBHelper()
880         print( db.db_selectKeyword() )
881         print( db.db_insertCrawlingData('1','2','3','4','5') )
882         db.db_free()
```