

```
1 -Pandas는 label이 부여된 data를 쉽고 직관적으로 취급할 수 있도록 설계된 Python third-party package이다.
2 -Pandas의 2가지 주요 data 구조인 Series(1차원 data)와 DataFrame(2차원 data)은 금융, 통계, 사회과학 등 많은 분야의 data 처리에 적합하다.
3 1. Pandas의 특징
4   1)pandas의 주요 기능을 설명
5     -쉬운 결손값(missing data)처리
6     -Label 위치를 자동적/명시적으로 정리한 data 작성
7     -Data 집약
8     -고도의 label base의 slicing, 추출, 큰 dataset의 subset화
9     -직감적인 dataset 결합
10    -Dataset의 유연한 변환 및 변형
11    -축의 계층적 label 붙임
12    -여러가지 data 형식에 대응한 강력한 I/O
13    -시계열 data 고유의 처리
14
15   2)pandas의 package 정보
16     -Version : 0.23.0
17     -공식 site : http://pandas.pydata.org
18     -Repository : https://github.com/pandas-dev/pandas
19     -PyPI : https://pypi.python.org/pypi/pandas
20
21   3)설치 여부 확인
22     !conda list | grep pandas    #Windows에서는 grep명령어를 사용할 수 없음.
23
24     -Windows에서는
25     !conda list    #목록 중에서 찾아야 함.
26
27
28 2. Series
29   1)Series는 index라고 불리는 label을 가진 동일한 data형을 가지는 1차원 data이다.
30   2)다음의 특징이 있다.
31     -Index(label)를 가지는 1차원 data
32     -Index는 중복 가능
33     -Label 또는 data의 위치를 지정한 추출가능. Index에 대한 slice가 가능
34     -산술 연산이 가능. 통계량을 산출하는 merit를 가지고 있음.
35
36   3)Python 표준 list나 tuple 등에서 사용되는 index라는 언어와의 혼동을 피하기 위해 series의 index를 label이라고 한다.
37
38   4)Series 작성하기
39     -Series의 작성에는 pandas.series class를 사용한다.
40     -제1인수에는 다음과 같은 1차원의 data를 넘겨준다.
41       --List
42       --Tuple
43       --Dirctionary
44       --numpy.ndarray
45     -아래와 같이 keyword 인수 index에 label이 되는 값을 넘기는 것으로 data를 표시한다.
46
47     import pandas as pd
48
49     ser = pd.Series([1,2,3], index=['a', 'b', 'c'])
```

```

50     ser
51     -----
52     a    1
53     b    2
54     c    3
55     dtype: int64
56
57     -index를 생략한 경우
58       --index를 생략한 경우에는 0부터 차례대로 정수가 할당된다.
59
60     ser = pd.Series([1,2,3])
61     ser
62     -----
63     0    1
64     1    2
65     2    3
66     dtype: int64
67
68     5)Label을 사용해서 data를 선택하기
69     -Series.loc를 사용해서 label에서 data를 선택한다.
70
71     ser.loc['b']
72     -----
73     2
74
75     -loc를 사용하지 않는 서식
76       --loc를 사용하지 않는 다음과 같은 서식도 있다.
77
78     ser['b']
79     -----
80     2
81
82     -Label의 범위 지정
83       --Label의 범위를 지정해서 slice를 할 수 있다.
84
85     ser.loc['b' : 'c']
86     -----
87     b    2
88     c    3
89     dtype: int64
90
91     --Lable에 따른 slice는 label의 시작 위치와 종료 위치를 포함한다.
92     --Python의 list나 tuple에 대한 slice와의 동작이 다른 경우에 주의한다.
93
94     -복수의 요소 지정
95       --복수의 요소를 list로 지정할 수 있다.
96
97     ser.loc[['a', 'c']]
98     -----
99     a    1
100    c    3
101    dtype: int64

```

```

102
103 6) 위치를 지정해서 data 선택하기
104   -Series.iloc를 사용해서 data의 위치를 정수값으로 지정하고 data를 선택할 수 있다.
105
106     ser.iloc[1]
107     -----
108     2
109
110   -iloc의 slice는 Python 표준 list나 tuple에 대한 slice와 같이 동작한다.
111   -위치를 slice로 지정
112
113     ser.iloc[1:3]
114     -----
115     b    2
116     c    3
117     dtype: int64
118
119 7) 논리값을 사용해서 data 선택하기
120   -loc와 iloc에는 논리값의 list를 넘길 수 있다.
121
122     ser.loc[[True, False, True]]
123     -----
124     a    1
125     c    3
126     dtype: int64
127
128   -인수에 부여된 논리값 list는 Series의 index 위치에 대하여 True에 지정된 위치만 되돌아간다.
129   -Series에 대한 비교 연산을 통해 논리값을 되돌려준다.
130
131     ser != 2
132     -----
133     a    True
134     b    False
135     c    True
136     dtype: bool
137
138   -이것을 이용해서 data를 추출할 수 있다.
139
140     ser.loc[ser != 2]
141     -----
142     a    1
143     c    3
144     dtype: int64
145
146
147 3. DataFrame
148 1) DataFrame은 행과 열에 label을 가진 2차원 data이다.
149 2) Data형은 열마다 다른 형을 가질 수 있다.
150 3) 1차원 data인 Series의 집합으로 인식하는 것도 가능하다.
151 4) Series의 특징을 포함해서 DataFrame에는 다음과 같은 특징이 있다.
152   -행과 열에 label을 가진 2차원 data
153   -열마다 다른 형태를 가질 수 있음.

```

```

154 -Table형 data에 대해 불러오기. data 쓰기가 가능
155 -DataFrame끼리 여러가지 조건을 사용한 결합 처리가 가능
156 -Cross 집계가 가능
157 5)Python 표준 list나 tuple 등에서 사용되는 index라는 언어와의 혼동을 피하기 위해서
    DataFrame의 index를 label이라고 한다.
158
159 6)DataFrame 작성하기
160 -DataFrame의 작성에는 pandas.DataFrame class를 사용한다.
161 -제1인수에는 1차원 또는 2차원 data를 넘긴다.
162 -Keyword 인수 index(행) 및 columns(열)에 label이 되는 값을 넘기는 것으로 data를 표시한
    다.
163
164 import pandas as pd
165
166 df = pd.DataFrame(
167     [[1, 10, 100], [2, 20, 200], [3, 30, 300]],
168     index=['r1', 'r2', 'r3'],
169     columns=['c1','c2','c3'])
170 df
171 -----
172      c1    c2    c3
173 r1     1    10  100
174 r2     2    20  200
175 r3     3    30  300
176
177 7)Label을 사용해서 data 선택하기
178 -Series와 같이 DataFrame.loc를 사용해서 data를 추출한다.
179 -DataFrame의 경우에는 행과 열의 label을 각각 지정한다.
180
181 df.loc['r2', 'c2']
182 -----
183 20
184
185 -모든 행(열)을 지정하는 경우
186 --모든 열을 지정하는 경우 요소에 [:]를 넘긴다.
187
188 df.loc['r2', :]
189 -----
190 c1     2
191 c2    20
192 c3   200
193 Name: r2, dtype: int64
194
195 --모든 행을 지정하는 경우도 같다.
196
197 df.loc[:, 'c2']
198 -----
199 r1    10
200 r2    20
201 r3    30
202 Name: c2, dtype: int64
203

```

```

204 -행의 data의 수가 1이고, 열의 data의 수가 복수 또는 행의 data의 수가 복수이고, 열의 data의
    수가 1의 경우. 되돌아오는 data형은 Series가 된다.
205 -Slice나 list를 넘겨주는 방법
206 --Slice나 list를 넘기는 방법은 Series와 같다.
207
208 # 행 label을 list로 지정, 열 label을 slice로 지정
209 df.loc[['r1', 'r3'], 'c2' : 'c3']
210 -----
211      c2    c3
212 r1    10  100
213 r3    30  300
214
215 -행의 data의 수와 열의 data의 수가 복수가 될 경우, 되돌아오는 data형은 DataFrame이다.
216
217 8)iloc를 사용해서 data를 선택하기
218 -Series와 같이 DataFrame.iloc를 사용해서 data의 위치에 의한 data 추출이 가능하다.
219 -DataFrame의 경우에는 행과 열의 위치를 각각 지정한다.
220
221 df.iloc[1:3, [0, 2]]
222 -----
223      c1    c3
224 r2     2  200
225 r3     3  300
226
227 9)열 이름을 지정해서 data 선택하기
228 -loc나 iloc를 지정하지 않고 DataFrame에 대해 아래와 같이 지정한 경우는 열 지정이 되고 되돌
    아 오는 Data형은 Series가 된다.
229
230 df['c2']
231 -----
232 r1    10
233 r2    20
234 r3    30
235 Name: c2, dtype: int64
236
237 10)논리값을 사용해서 data 선택하기
238 -Series와 같이 비교 연산을 하면 논리값이 되돌아온다.
239
240 df > 10
241 -----
242      c1      c2      c3
243 r1   False False  True
244 r2   False  True  True
245 r3   False  True  True
246
247 -DataFrame에서 Series를 사용해 비교 연산을 하여 data를 추출할 수 있다.
248
249 # c2열의 값이10보다 큰 data
250 df.loc[df['c2'] > 10]
251 -----
252      c1    c2    c3
253 r2     2   20  200

```

```

254     r3     3    30   300
255
256 -2개 이상의 조건을 조합하는 경우
257   --다음의 연산자를 사용한다.
258   --& : and 조건
259   --| : or 조건
260
261 -2개 이상의 조건을 조합하는 경우, 괄호()로 묶어 조건을 분류한다.
262
263   # c1열이 1보다 큰 동시에 c3열이 300보다 작은 data
264   df.loc[(df['c1'] > 1) & (df['c3'] < 300)]
265   -----
266         c1     c2     c3
267   r2      2     20    200
268
269 11)Series를 이용한 DataFrame 생성하기
270
271   list = [1,2,3]
272
273   ser1 = pd.core.series.Series(list)
274
275   ser2 = pd.core.series.Series(['one', 'two', 'three'])
276
277   pd.DataFrame(data=dict(num=ser1, word=ser2))
278   -----
279         num  word
280   0 1     one
281   1 2     two
282   2 3    three
283
284 12)Create DataFrame from your python code
285   -Python Dictionary를 이용해서 DataFrame 생성하기
286
287   emp_list = [
288       {'name':'John', 'age' : 25, 'job' : 'Manager'},
289       {'name':'Smith', 'age': 30, 'job' : 'Salesman'}
290   ]
291
292   df = pd.DataFrame(emp_list)
293   df
294   -----
295         age  job      name
296   0 25   Manager    John
297   1 30   Salesman  Smith      #list의 순서와 맞지 않음. key의 alphabet 순서
                                와 동일
298
299   df = df[['name','age','job']]
300
301   df.head()
302   -----
303         name  age  job
304   0   John   25  Manager

```

```

305     1   Smith   30   Salesman
306
307 -OrderDictionary 이용하기(key 순서 보장)
308
309     from collections import OrderedDict
310
311     emp_ordered_list = OrderedDict(
312         [
313             ('name', ['John', 'Smith']),
314             ('age', [25, 30]),
315             ('job', ['Manager', 'Salesman']),
316         ]
317     )
318
319     df = pd.DataFrame.from_dict(emp_ordered_list)
320
321     df.head()
322     -----
323         name   age   job
324     0 John     25   Manager
325     1 Smith    30   Salesman
326
327 -Python List를 이용해서 DataFrame 생성하기
328
329     emp_list = [
330         ['John', 25, 'Manager'],
331         ['Smith', 30, 'Salesman']
332     ]
333
334     column_names = ['name', 'age', 'job']
335
336     df = pd.DataFrame.from_records(emp_list, columns =column_names)
337
338     df.head()
339     -----
340         name   age   job
341     0 John     25   Manager
342     1 Smith    30   Salesman
343
344     emp_list = [
345         ['name', ['John', 'Smith']],
346         ['age', [25, 30]],
347         ['job', ['Manager', 'Salesman']],
348     ]
349
350     df = pd.DataFrame.from_items(emp_list)
351     -----
352     C:\ProgramData\Anaconda3\lib\site-packages\ipykernel_launcher.py:1:
FutureWarning: from_items is deprecated. Please use
DataFrame.from_dict(dict(items), ...) instead.
DataFrame.from_dict(OrderedDict(items)) may be used to preserve the key
order.

```

```
353 """Entry point for launching an IPython kernel.
```

```
354
```

```
355 df.head()
```

```
356
```

```
357     name  age  job
```

```
358 0 John    25  Manager
```

```
359 1 Smith   30  Salesman
```

```
360
```

```
361 13)DataFrame에서 column condition으로 select 하기
```

```
362
```

```
363 user_list = [
```

```
364     {'Name':'John', 'Age':25, 'Gender' : 'male', 'Address':'Chicago'},
```

```
365     {'Name':'Smith', 'Age':35, 'Gender' : 'male', 'Address':'Boston'},
```

```
366     {'Name':'Jenny', 'Age':45, 'Gender' : 'female', 'Address':'Dallas'}]
```

```
367
```

```
368
```

```
369 df = pd.DataFrame(user_list)
```

```
370
```

```
371 df = df[['Name', 'Age', 'Gender','Address']]
```

```
372
```

```
373 df.head()
```

```
374
```

```
375     Name  Age  Gender  Address
```

```
376 0 John    25   male   Chicago
```

```
377 1 Smith   35   male   'Boston'
```

```
378 2 Jenny   45  female   Dallas
```

```
379
```

```
380 -age가 30보다 많은 사람의 정보
```

```
381
```

```
382 df[df.Age > 30]
```

```
383
```

```
384     Name  Age  Gender  Address
```

```
385 1 Smith   35   male   Boston
```

```
386 2 Jenny   45  female   Dallas
```

```
387
```

```
388 df.query('Age > 30')
```

```
389
```

```
390     Name  Age  Gender  Address
```

```
391 1 Smith   35   male   Boston
```

```
392 2 Jenny   45  female   Dallas
```

```
393
```

```
394 -age가 30보다 많고 이름이 'Smith'인 사람의 정보
```

```
395
```

```
396 df[(df.Age > 30) & (df.Name == 'Smith')]
```

```
397
```

```
398     Name  Age  Gender  Address
```

```
399 1 Smith   35   male   Boston
```

```
400
```

```
401 14)Column filter 하기
```

```
402
```

```
403 -Column name이 Name과 Gender인 column만 가져오기
```

```
404
```



```

405     df.filter(items=['Name', 'Gender'])
406     -----
407         Name  Gender
408     0   John   male
409     1   Smith  male
410     2   Jenny female
411
412 -Column name이 'A'라는 글자가 있는 Column만 가져오기
413
414     df.filter(like = 'A')
415     -----
416         Age  Address
417     0    25   Chicago
418     1    35    Boston
419     2    45    Dallas
420
421 -정규식을 이용하여 column name이 'e'로 끝나는 Column만 가져오기
422
423     df.filter(regex = 'e$')
424     -----
425         Name  Age
426     0   John   25
427     1   Smith  35
428     2   Jenny  45
429
430 -정규식을 이용하여 column name이 'A'로 시작하는 Column만 가져오기
431
432     df.filter(regex = 'A^')
433
434 15)DataFrame의 row와 column 삭제하기
435 -row index를 사용하여 삭제하기
436
437     user_list = [
438         {'Age':25, 'Gender' : 'male', 'Address':'Chicago'},
439         {'Age':35, 'Gender' : 'male', 'Address':'Boston'},
440         {'Age':45, 'Gender' : 'female', 'Address':'Dallas'}
441     ]
442
443     df = pd.DataFrame(user_list,
444                       index = ['John', 'Smith', 'Jenny'],
445                       columns = ['Age', 'Gender', 'Address'])
446
447     df.head()
448     -----
449         Age  Gender  Address
450     John   25   male   Chicago
451     Smith 35   male    Boston
452     Jenny 45  female    Dallas
453
454     df.drop(['John', 'Jenny'])
455     -----
456         Age  Gender  Address

```

```

457     Smith  35   male      Boston
458
459     df.head()  #하지만 여전히 df는 예전값을 갖고 있다.
460
461     df = df.drop(['John', 'Jenny'])  #이렇게 하면 제거된 결과를 df가 갖게 된다.
462     df.head()
463     -----
464           Age  Gender  Address
465     Smith  35   male      Boston
466
467     -inplace keyword 인수 이용하기
468
469     df = pd.DataFrame(user_list,
470                       index = ['John', 'Smith', 'Jenny'],
471                       columns = ['Age', 'Gender', 'Address'])
472
473     df.drop(['John', 'Jenny'], inplace=True)
474     df
475     -----
476           Age  Gender  Address
477     Smith  35   male      Boston
478
479     -row가 숫자로 indexing되어 있을 때 삭제하기
480
481     user_list = [
482         {'Name':'John', 'Age':25, 'Gender' : 'male', 'Address':'Chicago'},
483         {'Name':'Smith', 'Age':35, 'Gender' : 'male', 'Address':'Boston'},
484         {'Name':'Jenny', 'Age':45, 'Gender' : 'female', 'Address':'Dallas'}
485     ]
486
487     df = pd.DataFrame(user_list)
488
489     df = df[['Name', 'Age', 'Gender', 'Address']]
490
491     df
492     -----
493           Name  Age  Gender  Address
494     0   John    25   male    Chicago
495     1   Smith   35   male    Boston
496     2   Jenny   45  female    Dallas
497
498     df = df.drop(df.index[[0, 2]])
499     df
500     -----
501           Name  Age  Gender  Address
502     1   Smith   35   male    Boston
503
504     -Condition으로 삭제하기
505
506     df.head()
507     -----
508           Name  Age  Gender  Address

```

```

509      0      John      25      male      Chicago
510      1      Smith      35      male      Boston
511      2      Jenny      45      female     Dallas

```

-Age가 30 이상인 사람의 정보만 저장함으로 나머지 정보는 삭제하는 방법

```

515      df = df[df.Age > 30]
516      df
517      -----
518           Name  Age  Gender  Address
519      1  Smith   35   male   Boston
520      2  Jenny   45  female   Dallas

```

-특정 column 삭제하기

```

524      df.head()
525      -----
526           Name  Age  Gender  Address
527      0  John    25   male   Chicago
528      1  Smith   35   male   Boston
529      2  Jenny   45  female   Dallas

```

531 df.drop('Age', axis= 1) <--- axis의 값이 1이면 column을 지칭

```

532      -----
533           Name  Gender  Address
534      0  John    male   Chicago
535      1  Smith   male   Boston
536      2  Jenny   female  Dallas

```

```

538      df.drop('Age', axis= 1, inplace=True)
539      df

```

```

540      -----
541           Name  Gender  Address
542      0  John    male   Chicago
543      1  Smith   male   Boston
544      2  Jenny   female  Dallas

```

546 16)DataFrame의 row, column 생성 및 수정하기

547 -column 새로 추가하기

```

549      df.head()
550      -----
551           Name  Age  Gender  Address Job
552      0  John    15   male   Chicago Student
553      1  Smith   25   male   Boston  Teacher
554      2  Jenny   17  female  Dallas  Student

```

```

556      df['Salary'] = 0

```

```

558      df.head()
559      -----
560           Name  Age  Gender  Address  Job      Salary

```

```

561      0   John    15   male      Chicago  Student    0
562      1   Smith   25   male      Boston   Teacher    0
563      2   Jenny   17  female     Dallas   Student    0

```

565 -기존 column 값을 이용

566 -'Job'에 따라 'Salary' 여부 Column으로 수정

```

567
568 df['Salary'] = np.where(df['Job'] != 'Student', 'yes', 'no')
569 df.head()

```

```

570 -----
571      Name  Age  Gender  Address Job      Salary
572  0   John   15   male      Chicago Student no
573  1   Smith  25   male      Boston   Teacher yes
574  2   Jenny  17  female     Dallas   Student no

```

576 -'Total' column 추가

```

577
578 student_list = [
579     {'Name':'John', 'Midterm':95, 'Final' : 85},
580     {'Name':'Smith', 'Midterm':85, 'Final' : 80},
581     {'Name':'Jenny', 'Midterm':30, 'Final' : 10},
582 ]

```

```

583
584 df = pd.DataFrame(student_list, columns= ['Name', 'Midterm', 'Final'])

```

```

585
586 df.head()

```

```

587 -----
588      Name  Midterm Final
589  0 John      95      85
590  1 Smith    85      80
591  2 Jenny    30      10

```

```

592
593 df['Total'] = df['Midterm'] + df['Final']
594 df

```

```

595 -----
596      Name  Midterm Final    Total
597  0 John      95      85     180
598  1 Smith    85      80     165
599  2 Jenny    30      10      40

```

601 -'Average' column 추가하기

```

602
603 df['Average'] = df['Total'] / 2
604 df.head()

```

```

605 -----
606      Name  Midterm Final    Total  Average
607  0 John      95      85     180     90.0
608  1 Smith    85      80     165     82.5
609  2 Jenny    30      10      40     20.0

```

611 -'Grade' column 추가하기

612

```

613     grade_list = []
614     for row in df['Average']:
615         if row <= 100 and row >= 90 :
616             grade_list.append('A')
617         elif row < 90 and row >= 80 :
618             grade_list.append('B')
619         elif row < 80 and row >= 70 :
620             grade_list.append('C')
621         elif row < 70 and row >= 60 :
622             grade_list.append('D')
623         else : grade_list.append('F')
624
625     df['Grade'] = grade_list
626     df
627     -----
628         Name  MidtermFinal      Total  AverageGrade
629     0 John      95      85      180      90.0      A
630     1 Smith    85      80      165      82.5      B
631     2 Jenny    30      10      40      20.0      F
632
633     -apply function 사용하기
634     -'Result' column 추가하기
635
636     def pass_or_fail(row):
637         if row != 'F':
638             return 'Pass'
639         else :
640             return 'Fail'
641
642     df['Result'] = df.Grade.apply(pass_or_fail)
643     df
644     -----
645         Name  MidtermFinal      Total  AverageGrade  Result
646     0 John      95      85      180      90.0      A      Pass
647     1 Smith    85      80      165      82.5      B      Pass
648     2 Jenny    30      10      40      20.0      F      Fail
649
650     -Column 추가하면서 각각의 값 조작하기
651
652     date_list = [
653         { 'yyyy-mm-dd' : '2019-01-05'},
654         { 'yyyy-mm-dd' : '2019-01-10'}
655     ]
656
657     df = pd.DataFrame(date_list, columns = ['yyyy-mm-dd'])
658     df
659     -----
660         yyyy-mm-dd
661     0 2019-01-05
662     1 2019-01-10
663
664     def extract_year(row):

```

```

665         return row.split('-')[0]
666
667     df['Year'] = df['yyyy-mm-dd'].apply(extract_year)
668     df
669     -----
670         yyyy-mm-dd      Year
671     0  2019-01-05    2019
672     1  2019-01-10    2019
673
674 -Row 추가하기
675
676     user_list = [
677         {'Name':'John', 'Age':15, 'Gender' : 'male', 'Address':'Chicago'},
678         {'Name':'Smith', 'Age':25, 'Gender' : 'male', 'Address':'Boston'},
679         {'Name':'Jenny', 'Age':17, 'Gender' : 'female', 'Address':'Dallas'}
680     ]
681
682     df = pd.DataFrame(user_list, columns = ['Name', 'Age', 'Gender', 'Address'])
683     df.head()
684     -----
685         Name  Age  Gender  Address
686     0  John   15   male   Chicago
687     1  Smith  25   male   Boston
688     2  Jenny  17  female   Dallas
689
690     df2 = pd.DataFrame(
691         [['Peter', 35, 'male', 'Seoul']], columns = ['Name', 'Age', 'Gender',
692         'Address']
693     )
694     df2.head()
695     -----
696         Name  Age  Gender  Address
697     0  Peter  35   male   Seoul
698
699     df.append(df2, ignore_index = True)
700     -----
701         Name  Age  Gender  Address
702     0  John   15   male   Chicago
703     1  Smith  25   male   Boston
704     2  Jenny  17  female   Dallas
705     3  Peter  35   male   Seoul
706
707
708 -passing keyword parameter to apply function
709
710     date_list = [{'Jumin': '2000-06-27'},
711         {'Jumin': '2002-09-24'},
712         {'Jumin': '2005-12-20'}]
713     df = pd.DataFrame(date_list, columns = ['Jumin'])
714     df
715     -----

```

```

716     Jumin
717     0 2000-06-27
718     1 2002-09-24
719     2 2005-12-20
720
721     def extract_year(row):
722         return row.split('-')[0]
723
724     df['Born_Year'] = df['Jumin'].apply(extract_year)
725     df
726     -----
727         Jumin      Born_Year
728     0 2000-06-27      2000
729     1 2002-09-24      2002
730     2 2005-12-20      2005
731
732     def calc_age(year, current_year):
733         return current_year - int(year)
734
735     df['Age'] = df['Born_Year'].apply(calc_age, current_year=2019)
736     df
737     -----
738         Jumin      Born_Year      Age
739     0 2000-06-27      2000         19
740     1 2002-09-24      2002         17
741     2 2005-12-20      2005         14
742
743     def get_introduce(age, prefix, suffix):
744         return prefix + str(age) + suffix
745
746     df['introduce'] = df['age'].apply(get_introduce, prefix="I am ", suffix=" years
747     old.")
748     df
749     -----
750         Jumin      Born_Year      Age      Introduce
751     0 2000-06-27      2000         19      I am 19 years old.
752     1 2002-09-24      2002         17      I am 17 years old.
753     2 2005-12-20      2005         14      I am 14 years old.
754
755     17)How to use map function
756     -apply()와 유사한 map()
757
758     date_list = [{'Date': '2000-06-27'},
759                  {'Date': '2002-09-24'},
760                  {'Date': '2005-12-20'}]
761     df = pd.DataFrame(date_list, columns = ['Date'])
762     df
763     -----
764         Date
765     0 2000-06-27
766     1 2002-09-24
767     2 2005-12-20

```

```

767
768     def extract_year(date):
769         return date.split('-')[0]
770
771     df['Year'] = df['Date'].map(extract_year)
772     df
773     -----
774         Date      Year
775     0 2000-06-27   2000
776     1 2002-09-24   2002
777     2 2005-12-20   2005
778
779 -map() 응용하기
780
781     data_list = [
782         {'Name':'John', 'Age':15, 'Gender':'male', 'Job':'Student'},
783         {'Name':'Smith', 'Age':25, 'Gender':'male', 'Job':'Teacher'},
784         {'Name':'Jenny', 'Age':27, 'Gender':'female', 'Job':'Developer'},
785     ]
786
787     df = pd.DataFrame(data_list, columns = ['Name', 'Age', 'Gender', 'Job'])
788     df
789     -----
790         Name  Age  Gender Job
791     0 John    15   male   Student
792     1 Smith   25   male   Teacher
793     2 Jenny   27  female  Developer
794
795     df.Job = df.Job.map({'Student':1, 'Teacher':2, 'Developer':3})
796     df
797     -----
798         Name  Age  Gender Job
799     0 John    15   male    1
800     1 Smith   25   male    2
801     2 Jenny   27  female    3
802
803 18)Applymap 사용하기
804     -한번에 DataFrame에 있는 모든 요소들을 수정하고자 할 때
805
806     data_list = [
807         {'x': 5.5, 'y': -5.6},
808         {'x': -5.2, 'y': 5.5},
809         {'x': -1.6, 'y': -4.5}
810     ]
811     df = pd.DataFrame(data_list)
812     df
813     -----
814         x      y
815     0  5.5  -5.6
816     1 -5.2   5.5
817     2 -1.6  -4.5
818

```



```

819     import numpy as np
820
821     df = df.applymap(np.around) #반올림함수 사용
822     df
823     -----
824         x      y
825     0   6.0  -6.0
826     1  -5.0   6.0
827     2  -2.0  -4.0
828
829 19)Group By
830
831 student_list = [
832     {'Name': 'John', 'Major': "Computer Science", 'Gender': "male"},
833     {'Name': 'Nate', 'Major': "Computer Science", 'Gender': "male"},
834     {'Name': 'Abraham', 'Major': "Physics", 'Gender': "male"},
835     {'Name': 'Brian', 'Major': "Psychology", 'Gender': "male"},
836     {'Name': 'Janny', 'Major': "Economics", 'Gender': "female"},
837     {'Name': 'Yuna', 'Major': "Economics", 'Gender': "female"},
838     {'Name': 'Jeniffer', 'Major': "Computer Science", 'Gender': "female"},
839     {'Name': 'Edward', 'Major': "Computer Science", 'Gender': "male"},
840     {'Name': 'Zara', 'Major': "Psychology", 'Gender': "female"},
841     {'Name': 'Wendy', 'Major': "Economics", 'Gender': "female"},
842     {'Name': 'Sera', 'Major': "Psychology", 'Gender': "female"}
843 ]
844 df = pd.DataFrame(student_list, columns = ['Name', 'Major', 'Gender'])
845 df
846 -----
847      Name  Major      Gender
848 0   John  Computer Science  male
849 1   Nate  Computer Science  male
850 2 Abraham  Physics        male
851 3   Brian  Psychology      male
852 4   Janny  Economics      female
853 5   Yuna  Economics      female
854 6 Jeniffer Computer Science  female
855 7   Edward Computer Science  male
856 8    Zara  Psychology      female
857 9   Wendy  Economics      female
858 10  Sera  Psychology      female
859
860 -전공별 Group By
861     groupby_major = df.groupby('Major')
862
863     groupby_major.groups
864     -----
865     {'Computer Science': Int64Index([0, 1, 6, 7], dtype='int64'),
866     'Economics': Int64Index([4, 5, 9], dtype='int64'),
867     'Physics': Int64Index([2], dtype='int64'),
868     'Psychology': Int64Index([3, 8, 10], dtype='int64')}
869
870 -보기 좋게

```

```

871     for name, group in groupby_major:
872         print(name + ": " + str(len(group)))
873         print(group)
874         print()
875     -----
876     Computer Science: 4
877         Name    Major    Gender
878     0    John    Computer Science    male
879     1    Nate    Computer Science    male
880     6    Jeniffer    Computer Science    female
881     7    Edward    Computer Science    male
882
883     Economics: 3
884         Name    Major    Gender
885     4    Janny    Economics    female
886     5    Yuna    Economics    female
887     9    Wendy    Economics    female
888
889     Physics: 1
890         Name    Major    Gender
891     2    Abraham    Physics    male
892
893     Psychology: 3
894         Name    Major    Gender
895     3    Brian    Psychology    male
896     8    Zara    Psychology    female
897     10    Sera    Psychology    female
898
899     -전공별 명수
900
901     df_major_cnt = pd.DataFrame({'Count':groupby_major.size()})
902     df_major_cnt
903     -----
904         Count
905     Major
906     Computer Science    4
907     Economics          3
908     Physics             1
909     Psychology          3
910
911     -전공도 count와 같이
912
913     df_major_cnt = pd.DataFrame({'Count':groupby_major.size()}).reset_index()
914     df_major_cnt
915     -----
916         Major    Count
917     0    Computer Science    4
918     1    Economics          3
919     2    Physics             1
920     3    Psychology          3
921
922     -성별로 group by

```

```

923
924     groupby_gender = df.groupby('Gender')
925
926     for name, group in groupby_gender:
927         print(name + ": " + str(len(group)))
928         print(group)
929         print()
930     -----
931     female: 6
932         Name      Major      Gender
933     4   Janny      Economics    female
934     5   Yuna       Economics    female
935     6   Jeniffer   Computer Science  female
936     8   Zara       Psychology    female
937     9   Wendy     Economics    female
938     10  Sera       Psychology    female
939
940     male: 5
941         Name      Major      Gender
942     0   John      Computer Science  male
943     1   Nate      Computer Science  male
944     2   Abraham   Physics          male
945     3   Brian     Psychology    male
946     7   Edward   Computer Science  male
947
948     20)중복된 값 제거하기
949
950     data_list = [
951         {'Name':'John', 'Gender':'male', 'Job':'Student'},
952         {'Name':'Smith','Gender':'male', 'Job':'Teacher'},
953         {'Name':'Jenny','Gender':'female', 'Job':'Developer'},
954         {'Name':'Smith','Gender':'male', 'Job':'Teacher'}
955     ]
956     df = pd.DataFrame(data_list, columns = ['Name', 'Gender', 'Job'])
957
958     df.head()
959     -----
960         Name Gender   Job
961     0   John   male   Student
962     1   Smith male   Teacher
963     2   Jenny female Developer
964     3   Smith male   Teacher      #중복된 값
965
966     -중복된 값 확인하기
967
968     df.duplicated()
969     -----
970     0   False
971     1   False
972     2   False
973     3    True      #여기가 중복된 값이 있다는 뜻
974     dtype: bool

```

```

975
976 -중복된 값 제거
977
978 df.drop_duplicates()
979 -----
980      Name Gender   Job
981 0 John    male   Student
982 1 Smith male    Teacher
983 2 Jenny female Developer
984
985 21)None Value(NaN)을 찾고 원하는 값으로 변경하기
986
987 student_list = [
988     {'Name': 'John', 'Major': 'Computer Science', 'Gender': 'male', 'Age': 40},
989     {'Name': 'Nate', 'Major': None, 'Gender': "male", 'Age':35},
990     {'Name': 'Abraham', 'Major': 'Physics', 'Gender': 'male', 'Age':37},
991     {'Name': 'Brian', 'Major': 'Psychology', 'Gender': 'male', 'Age':None},
992     {'Name': 'Janny', 'Major': None, 'Gender': 'female', 'Age':10},
993     {'Name': 'Yuna', 'Major': None, 'Gender': 'female', 'Age':12},
994     {'Name': 'Jeniffer', 'Major': 'Computer Science', 'Gender': 'female',
995     'Age':45},
996     {'Name': 'Edward', 'Major': 'Computer Science', 'Gender': 'male',
997     'Age':None},
998     {'Name': 'Zara', 'Major': 'Psychology', 'Gender': 'female', 'Age':25},
999     {'Name': 'Wendy', 'Major': 'Economics', 'Gender': 'female', 'Age':37},
1000     {'Name': 'Sera', 'Major': None, 'Gender': 'female', 'Age':None}
1001 ]
1002 df = pd.DataFrame(student_list, columns = ['Name', 'Major', 'Gender', 'Age'])
1003 df
1004 -----
1005      Name  Major      Gender  Age
1006 0  John  Computer Science  male    40.0
1007 1  Nate      None         male    35.0
1008 2 Abraham Physics         male    37.0
1009 3  Brian Psychology        male    NaN
1010 4  Janny  None          female    10.0
1011 5  Yuna   None          female    12.0
1012 6 Jeniffer Computer Science female    45.0
1013 7  Edward Computer Science  male     NaN
1014 8  Zara   Psychology        female    25.0
1015 9  Wendy  Economics         female    37.0
1016 10 Sera   None            female    NaN
1017
1018 df.shape
1019 -----
1020 (11,4)
1021
1022 df.info()
1023 -----
1024 <class 'pandas.core.frame.DataFrame'>
RangeIndex: 11 entries, 0 to 10
Data columns (total 4 columns):

```

```

1025 Name    11 non-null object
1026 Major    7 non-null object
1027 Gender   11 non-null object
1028 Age      8 non-null float64
1029 dtypes: float64(1), object(3)
1030 memory usage: 432.0+ bytes
1031
1032 -Null 확인하기
1033
1034 df.isna()
1035 -----
1036      Name Major Gender Age
1037 0  False False  False  False
1038 1  False  True   False  False
1039 2  False False  False  False
1040 3  False False  False   True
1041 4  False  True   False  False
1042 5  False  True   False  False
1043 6  False False  False  False
1044 7  False False  False   True
1045 8  False False  False  False
1046 9  False False  False  False
1047 10 False  True   False   True
1048
1049 df.isnull()
1050 -----
1051      Name Major Gender Age
1052 0  False False  False  False
1053 1  False  True   False  False
1054 2  False False  False  False
1055 3  False False  False   True
1056 4  False  True   False  False
1057 5  False  True   False  False
1058 6  False False  False  False
1059 7  False False  False   True
1060 8  False False  False  False
1061 9  False False  False  False
1062 10 False  True   False   True
1063
1064 -None값을 다른 값으로 변경하기
1065
1066 df.Age = df.Age.fillna(0)    #숫자 NaN을 0으로
1067 df.Major = df.Major.fillna('Unknown') #글자 None을 Unknown으로
1068 df
1069 -----
1070      Name  Major      Gender  Age
1071 0  John  Computer Science  male    40.0
1072 1  Nate    Unknown      male    35.0
1073 2 Abraham Physics      male    37.0
1074 3  Brian Psychology    male     0.0
1075 4  Janny   Unknown    female    10.0
1076 5  Yuna    Unknown    female    12.0

```

```

1077      6  Jeniffer Computer Science  female    45.0
1078      7   Edward Computer Science   male      0.0
1079      8    Zara    Psychology        female    25.0
1080      9   Wendy    Economics        female    37.0
1081     10   Sera      Unknown          female     0.0

```

22) Unique와 갯수 알아보기

-Unique 즉, 중복제거된 값 알아보기

```
df.Major.unique()
```

```
-----
array(['Computer Science', 'Unknown', 'Physics', 'Psychology', 'Economics'],
      dtype=object)
```

```
df.Gender.unique()
```

```
-----
array(['male', 'female'], dtype=object)
```

```
df.Major.value_counts()
```

```
-----
Unknown          4
Computer Science  3
Psychology        2
Physics           1
Economics         1
Name: Major, dtype: int64
```

23) DataFrame 합치기

```
list1 = [
    {'Name': 'John', 'Job': 'Teacher'},
    {'Name': 'Nate', 'Job': 'Student'},
    {'Name': 'Fred', 'Job': 'Developer'}
]
```

```
list2 = [
    {'Name': 'Ed', 'Job': 'Dentist'},
    {'Name': 'Jack', 'Job': 'Farmer'},
    {'Name': 'Ted', 'Job': 'Designer'}
]
```

```
df1 = pd.DataFrame(list1, columns = ['Name', 'Job'])
```

```
df2 = pd.DataFrame(list2, columns = ['Name', 'Job'])
```

-concat()로 합치기

```
result = pd.concat([df1, df2])
```

```
result
```

```
-----
   Name Job
0  John  Teacher
1  Nate  Student
```

```

1128     2   Fred   Developer
1129     0   Ed    Dentist          #0, 1, 2가 반복
1130     1   Jack   Farmer
1131     2   Ted    Designer
1132
1133     result = pd.concat([df1, df2], ignore_index = True)
1134     result
1135     -----
1136           Name Job
1137     0   John   Teacher
1138     1   Nate   Student
1139     2   Fred   Developer
1140     3   Ed    Dentist
1141     4   Jack   Farmer
1142     5   Ted    Designer
1143
1144     -append()로 합치기
1145
1146     result = df1.append(df2)
1147     result
1148     -----
1149           Name Job
1150     0   John   Teacher
1151     1   Nate   Student
1152     2   Fred   Developer
1153     0   Ed    Dentist          #0, 1, 2가 반복
1154     1   Jack   Farmer
1155     2   Ted    Designer
1156
1157     result = pd.append([df1, df2], ignore_index = True)
1158     result
1159     -----
1160           Name Job
1161     0   John   Teacher
1162     1   Nate   Student
1163     2   Fred   Developer
1164     3   Ed    Dentist
1165     4   Jack   Farmer
1166     5   Ted    Designer
1167
1168     -Column으로 합치기
1169     --두개의 DataFrame의 column이 서로 일치하지 않음.
1170
1171     list1 = [
1172         {'Name': 'John', 'Job': 'Teacher'},
1173         {'Name': 'Nate', 'Job': 'Student'},
1174         {'Name': 'Jack', 'Job': 'Developer'}
1175     ]
1176
1177     list2 = [
1178         {'Age': 25, 'Country': 'U.S'},
1179         {'Age': 30, 'Country': 'U.K'},

```

```

1180         {'Age': 45, 'Country': 'Korea'}
1181     ]
1182
1183     df1 = pd.DataFrame(list1, columns = ['Name', 'Job'])
1184     df2 = pd.DataFrame(list2, columns = ['Age', 'Country'])
1185
1186     result = pd.concat([df1, df2], axis=1, ignore_index=True)
1187     result
1188     -----
1189           0      1      2      3
1190    0   John  Teacher   25   U.S
1191    1    Nate   Student   30   U.K
1192    2   Jack   Developer  45   Korea

```

1195 4. 다양한 data 불러오기

```

1196     1)Pandas는 다음과 같이 다양한 형식의 data를 불러올 수 있다.
1197     -CSV
1198     -Excel
1199     -Database
1200     -JSON
1201     -MessagePack
1202     -HTML
1203     -Google BigQuery
1204     -Clipboard
1205     -Pickle
1206     -기타(http://pandas.pydata.org/pandas-docs/stable/io.html)

```

1209 5. CSV file 불러오기

```

1210     1)pandas.read_csv() 함수를 사용한다.
1211     2)첫번째 인수에 file 경로를 넘겨주면 DataFrame 형 object를 넘겨준다.
1212     3)File 경로 또는 URL 형식으로 지정할 수 있다.

```

```

1214     import pandas as pd
1215
1216     df = pd.read_csv('pandas_data/friend_list.csv')
1217     df.head()
1218     -----
1219           name  age  job
1220    0   John    20  student
1221    1  Jenny    30  developer
1222    2   Nate    30  teacher
1223    3  Julia    40  dentist
1224    4  Brian    45  manager

```

```

1226     4)DataFrame.head()는 앞에서 5행분의 DataFrame을 넘겨준다.
1227     5)인수에 정수값을 넘기는 방식으로 행수를 지정할 수 있다.

```

```

1229     df = pd.read_csv('pandas_data/friend_list.csv')
1230     df.head(2)
1231     -----

```



```

1232     name age  job
1233  0 John   20  student
1234  1 Jenny 30  developer
1235

```

6)뒤에서부터 읽을 행의 갯수를 지정할 수도 있다.

```

1237
1238 df = pd.read_csv('pandas_data/friend_list.csv')
1239 df.tail(2)

```

```

1240 -----
1241     name  age  job
1242  4  Brian  45  manager
1243  5  Chris  25  intern
1244

```

7)각 열은 Series이다.

```

1246
1247 type(df.job)
1248 -----
1249 pandas.core.series.Series
1250

```

8)지정한 열을 DataFrame의 index로 하기

-아래 예제처럼 keyword 인수 index_col에 수치 또는 열 이름을 지정하여 지정한 열을 DataFrame의 index로 한다.

```

1253
1254 # index로 지정할 열을 번호로 지정
1255 df = pd.read_csv('pandas_data/friend_list.csv', index_col = 0)
1256 df.head()

```

```

1257 -----
1258         age  job
1259  name
1260 John    20  student
1261 Jenny   30  developer
1262 Nate    30  teacher
1263 Julia   40  dentist
1264 Brian   45  manager
1265

```

```

1266 #index로 지정할 열을 열 이름으로 지정
1267 df = pd.read_csv('pandas_data/friend_list.csv', index_col = 'job')
1268 df.head()

```

```

1269 -----
1270         name  age
1271  job
1272 student    John    20
1273 developer  Jenny    30
1274 teacher   Nate     30
1275 dentist   Julia    40
1276 manager   Brian    45
1277

```

9)지정한 열을 지정한 형으로 불러오기

-Keyword 인수 dtype에 열 이름(key)과 형(값)을 사전형으로 지정하여 지정한 열을 지정한 형으로 불러올 수 있다.

```

1280
1281 #형 지정

```

```
1282 df = pd.read_csv('pandas_data/friend_list.csv', dtype={'age':float})
1283 df.head()
```

```
1284 -----
```

```
1285      name  age  job
1286  0   John  20.0  student
1287  1  Jenny  30.0  developer
1288  2   Nate  30.0  teacher
1289  3  Julia  40.0  dentist
1290  4  Brian  45.0  manager
```

```
1291
```

```
1292 10)형식이 유사한 txt file 읽어오기
```

```
1293 -CSV file처럼 txt file도 각 열의 구분을 ','로 할 경우
```

```
1294
```

```
1295 df = pd.read_csv('pandas_data/friend_list.txt')
1296 df.head()
```

```
1297 -----
```

```
1298      name  age  job
1299  0   John   20  student
1300  1  Jenny   30  developer
1301  2   Nate   30  teacher
1302  3  Julia   40  dentist
1303  4  Brian   45  manager
```

```
1304
```

```
1305
```

```
1306 11)만일 file의 column들이 쉼표로 구분되어 있지 않은 경우
```

```
1307 -delimiter parameter에 구분자를 지정해서 column을 나눠야 한다.
```

```
1308
```

```
1309 df = pd.read_csv('pandas_data/friend_list_tab.txt')
1310 df.head()
```

```
1311 -----
```

```
1312      name age job
1313 0 John\t20\tstudent
1314 1 Jenny\t30\tdeveloper
1315 2 Nate\t30\tteacher
1316 3 Julia\t40\tdentist
1317 4 Brian\t45\tmanager    #구분이 어려움.
```

```
1318
```

```
1319 df = pd.read_csv('pandas_data/friend_list_tab.txt', delimiter = '\t')
1320 df.head()
```

```
1321 -----
```

```
1322      name  age  job
1323  0   John   20  student
1324  1  Jenny   30  developer
1325  2   Nate   30  teacher
1326  3  Julia   40  dentist
1327  4  Brian   45  manager
```

```
1328
```

```
1329 12)file에 data header가 없는 csv file을 사용할 때
```

```
1330 -만일 data header가 없으면, header = None으로 지정해야 첫번째 data가 data header로
    들어가는 것을 막을 수 있다.
```

```
1331
```

```
1332 df = pd.read_csv('pandas_data/friend_list_no_head.csv')
```

```

1333 df.head()
1334 -----
1335      John    20  student      #첫 번째 행이 header가 돼버림.
1336 0    Jenny 30  developer
1337 1     Nate  30  teacher
1338 2     Julia 40  dentist
1339 3     Brian 45  manager
1340 4      Chris 25   intern
1341
1342 df = pd.read_csv('pandas_data/friend_list_no_head.csv', header = None)
1343 df
1344 -----
1345      0      1      2
1346 0    John    20  student
1347 1    Jenny 30  developer
1348 2     Nate  30  teacher
1349 3     Julia 40  dentist
1350 4     Brian 45  manager
1351 5      Chris 25   intern
1352
1353 -만일 header가 없는 data를 호출했을 경우, DataFrame 생성 후, column header를 지정할
    수 있다.
1354
1355 df.columns = ['Name', 'Age', 'Job']
1356 df.index = ['1101', '1102', '1103', '1104', '1105', '1106']
1357 df
1358 -----
1359      Name  Age  Job
1360 1101  John   20  student
1361 1102  Jenny  30  developer
1362 1103  Nate   30  teacher
1363 1104  Julia  40  dentist
1364 1105  Brian  45  manager
1365 1106  Chris  25   intern
1366
1367 -file을 열 때 동시에 header에 column을 지정해야 할 경우
1368
1369 df = pd.read_csv('pandas_data/friend_list_no_head.csv', header = None,
1370 names=['Name', 'Age', 'Job'])
1371 df.head()
1372 -----
1373 위의 결과와 동일
1374
1375 13)외부 CSV file 읽기
1376 -https://github.com/vincentarelbundock/Rdatasets/tree/master/csv/datasets
1377
1378 14)기타 option
1379 -read_csv()에는 다수의 option이 준비되어 있다.
1380 -상세한 문서의 내용은 문서를 참조한다.
1381 -http://panda.pydata.org/pandas-docs/stable/io.html#io-read-csv-table
1382

```

```

1383 15)DataFrame csv file로 저장하기
1384
1385 import pandas as pd
1386
1387 user_list = [
1388     {'Name':'John', 'Age':25, 'Gender' : 'male', 'Address':'Chicago'},
1389     {'Name':'Smith', 'Age':35, 'Gender' : 'male', 'Address':None},
1390     {'Name':'Jenny', 'Age':45, 'Gender' : 'female', 'Address':'Dallas'}
1391 ]
1392
1393 df = pd.DataFrame(user_list)
1394
1395 df = df[['Name', 'Age', 'Gender','Address']]
1396
1397 df.head()
1398 -----
1399 Name    Age  Gender  Address
1400 0 John    25   male    Chicago
1401 1 Smith   35   male      None
1402 2 Jenny   45  female    Dallas
1403
1404 df.to_csv('user_list.csv')
1405
1406 -기본적으로 to_csv()는 index = True, header = True가 설정되어 있다.
1407 --만일 index = False로 설정할 경우
1408
1409     df.to_csv('user_list.csv', index=False)
1410
1411 --각 행의 index 즉 0, 1, 2가 사라진 csv file이 생성된다.
1412 --또한 header = False로 설정하면 각 열의 header가 없어진다.
1413
1414 -Smith의 거주지가 None이기 때문에 빈칸으로 값이 들어간다.
1415 -만일 csv로 file 저장시 빈칸대신 '-'를 넣어서 뭔가 있다는 것을 설정하려면,
1416
1417     df.to_csv('user_list.csv', na_rep = '-')
1418
1419 -이렇게 설정하면 해당 칸에는 '-'가 들어가게 된다.
1420
1421
1422 6. Excel file 불러오기
1423 1)pandas.read_excel() 함수를 사용한다.
1424 2)사전준비를 위해 xlrd module을 설치여부를 확인한다.
1425
1426 !conda list | grep xlrd    #Windows에서는 사용 불가
1427
1428 $ conda install -y xlrd==1.2.0    #설치안되어 있으면 설치
1429
1430 # Excel file 불러오기
1431 df = pd.read_excel('pandas_data/재무실적.xlsx')
1432 df.head()
1433
1434 3)불러오는 sheet 지정하기

```

```

1435     -기본 설정으로는 첫 번째 sheet를 불러온다.
1436     -Sheet 이름을 지정해서 불러올 때는 keyword 인수 sheetname에 sheet이름을 지정한다.
1437
1438     df = pd.read_excel('pandas_data/재무실적.xlsx', sheetname = 'data2')
1439     df.head()
1440
1441
1442 7. SQL을 사용해서 불러오기
1443     1)pandas.read_sql() 함수를 사용한다.
1444     2)첫번째 인수에 query를 실행하는 SQL문, 두번째 인수에
1445     SQLAlchemy(http://docs.sqlalchemy.org/en/latest/dialects/index) 또는 DBAPI2(PEP
1446     249, Python Database API Specification v2.0,
1447     https://www.python.org/dev/peps/pep-0249)의 접속 instance를 넘겨준다.
1448
1449     3)MariaDB with Python
1450     -설치여부 확인하기
1451
1452     !conda list | grep mysql-connector-python      #Windows에서는 grep 명령어 사
1453     용 불가
1454
1455     -mysql-connector-python 설치하기
1456
1457     --In Anaconda Prompt,
1458
1459     $ conda install -y mysql-connector-python
1460
1461     import mysql.connector as mariadb
1462
1463     mariadb_connection = mariadb.connect(user='root',
1464     password='javamariadb', host='localhost', database='world')
1465     cursor = mariadb_connection.cursor()
1466
1467     cursor.execute("SELECT ID, Name, CountryCode, District, Population FROM
1468     city WHERE CountryCode='KOR'")
1469
1470     for ID,Name,CountryCode,District,Population in cursor:
1471     print('ID = %d, Name = %s, CountryCode = %s, District = %s,
1472     Popluation = %d' % (ID, Name, CountryCode,District,Population))
1473     -----
1474     ID = 2331, Name = Seoul, CountryCode = KOR, District = Seoul, Popluation
1475     = 9981619
1476     ID = 2332, Name = Pusan, CountryCode = KOR, District = Pusan,
1477     Popluation = 3804522
1478     ID = 2333, Name = Inchon, CountryCode = KOR, District = Inchon,
1479     Popluation = 2559424
1480     ID = 2334, Name = Taegu, CountryCode = KOR, District = Taegu,
1481     Popluation = 2548568
1482     ID = 2335, Name = Taejon, CountryCode = KOR, District = Taejon,
1483     Popluation = 1425835
1484     ...
1485     ...
1486     4)Oracle with Python

```

```

1475 -Oracle cx_Oracle 7
1476 -https://www.oracle.com/database/technologies/appdev/python.html
1477 -Install on Windows
1478   $ python -m pip install cx_Oracle --upgrade
1479
1480 -In Anaconda Prompt
1481   $ conda install cx_oracle
1482
1483 -Connection
1484   import cx_Oracle
1485
1486   --conn = cx_Oracle.connect('scott', 'tiger', 'localhost:1521/XE')
1487
1488   --conn1 = cx_Oracle.connect('scott/tiger@localhost:1521/XE')
1489
1490   print(conn1)
1491   -----
1492   <cx_Oracle.Connection to scott@localhost:1521/XE>
1493
1494   --dsn_tns = cx_Oracle.makedsn('localhost', 1521, 'XE')
1495   print(dsn_tns)
1496   -----
1497   (DESCRIPTION=(ADDRESS=(PROTOCOL=TCP)(HOST=localhost)(PORT=152
1498   1))(CONNECT_DATA=(SID=XE)))
1499
1500   conn2 = cx_Oracle.connect('scott', 'tiger', dsn_tns)
1501   print(conn2)
1502   -----
1503   <cx_Oracle.Connection to
1504   scott@\(DESCRIPTION=\(ADDRESS=\(PROTOCOL=TCP\)\(HOST=localhost\)\(POR
1505   T=1521\)\)\(CONNECT\_DATA=\(SID=XE\)\)\)>
1506
1507   conn.version
1508   -----
1509   '11.2.0.2.0'
1510
1511 -Cursor Objects
1512   cursor = conn.cursor()
1513   cursor.execute('SELECT empno, ename, job, hiredate, mgr, sal, comm,
1514   deptno FROM emp')
1515   -----
1516   <cx_Oracle.Cursor on <cx_Oracle.Connection to
1517   scott@localhost:1521/XE>>
1518
1519   for empno, ename, job, hiredate, mgr, sal, comm, deptno in cursor:
1520       print(empno, ename, job, hiredate, mgr, sal, comm,deptno)
1521   -----
1522   7369 SMITH CLERK 1980-12-17 00:00:00 7902 800.0 None 20
1523   7499 ALLEN SALESMAN 1981-02-20 00:00:00 7698 1600.0 300.0 30
1524   7521 WARD SALESMAN 1981-02-22 00:00:00 7698 1250.0 500.0 30
1525   7566 JONES MANAGER 1981-04-02 00:00:00 7839 2975.0 None 20
1526   7654 MARTIN SALESMAN 1981-09-28 00:00:00 7698 1250.0 1400.0 30

```

```

1522      7698 BLAKE MANAGER 1981-05-01 00:00:00 7839 2850.0 None 30
1523      7782 CLARK MANAGER 1981-06-09 00:00:00 7839 2450.0 None 10
1524      7788 SCOTT ANALYST 1987-07-13 00:00:00 7566 3000.0 None 20
1525      7839 KING PRESIDENT 1981-11-17 00:00:00 None 5000.0 None 10
1526      7844 TURNER SALESMAN 1981-09-08 00:00:00 7698 1500.0 0.0 30
1527      7876 ADAMS CLERK 1987-07-13 00:00:00 7788 1100.0 None 20
1528      7900 JAMES CLERK 1981-12-03 00:00:00 7698 950.0 None 30
1529      7902 FORD ANALYST 1981-12-03 00:00:00 7566 3000.0 None 20
1530      7934 MILLER CLERK 1982-01-23 00:00:00 7782 1300.0 None 10

```

```

1531
1532      cursor.close()

```

```

1533
1534

```

1535 8. HTML file 불러오기

1536 1)pandas.read_html() 함수 사용하기

1537 -다음의 third-party package들의 설치여부를 확인한다.

1538 --html5lib, lxml, BeautifulSoup4

1539 -HTML file의 table 요소를 DataFrame으로 불러온다.

1540 -DataFrame이 들어 있던 list가 return된다.

1541 -Table 요소가 여러 개 있는 경우, 여러 개의 DataFrame이 저장된다.

1542

```

1543      import pandas as pd

```

1544

```

1545      url = 'http://www.fdic.gov/bank/individual/failed/banklist.html'

```

1546

```

1547      dfs = pd.read_html(url)

```

1548

```

1549      dfs[0].info()

```

1550

```

--

```

```

1551      <class 'pandas.core.frame.DataFrame'>

```

```

1552      RangeIndex: 555 entries, 0 to 554

```

```

1553      Data columns (total 7 columns):

```

```

1554      Bank Name          555 non-null object

```

```

1555      City              555 non-null object

```

```

1556      ST               555 non-null object

```

```

1557      CERT            555 non-null int64

```

```

1558      Acquiring Institution 555 non-null object

```

```

1559      Closing Date      555 non-null object

```

```

1560      Updated Date      555 non-null object

```

```

1561      dtypes: int64(1), object(6)

```

```

1562      memory usage: 30.4+ KB

```

1563

1564 2)request library 이용하기

1565 -Python HTTP for Humans'

1566 -<http://docs.python-requests.org/en/master/>

1567 -Install

```

1568      pip install requests

```

1569

```

1570      import requests

```

```

1571      url = 'https://www.naver.com'

```

```

1572      naver = requests.get(url)

```

```

1573     print(naver.text)
1574
1575 3) pprint module 이용하기
1576     -대량의 data를 보기 쉽게 표시해주는 표준 module
1577     -pprint(prettyprint)
1578
1579     import requests
1580     import pprint
1581     url = 'https://www.naver.com'
1582     naver = requests.get(url)
1583     pprint.pprint(naver.text)
1584
1585 4)requests를 사용하여 API에 접근하기
1586     -기상청 RSS(http://www.weather.go.kr/weather/lifenindustry/sevice\_rss.jsp)를 이
1587     용하자.
1588     -2013년 동네 예보 RSS
1589     http://www.weather.go.kr/images/weather/lifenindustry/dongnaeforecast\_rss.
1590     pdf
1591     -주소선택후 rss button click하면 zone을 알 수 있다.
1592     --예:서울특별시 강남구 청담동
1593     http://www.kma.go.kr/wid/queryDFSRSS.jsp?zone=1168056500
1594
1595     import requests
1596     import pprint
1597     api_uri = 'http://www.kma.go.kr/wid/queryDFSRSS.jsp?zone=1168056500'
1598     weather_data = requests.get(api_uri).text
1599     pprint.pprint(weather_data)
1600
1601     -----
1602     ('<?xml version="1.0" encoding="UTF-8" ?>\n'
1603     '<rss version="2.0">\n'
1604     '<channel>\n'
1605     '<title>기상청 동네예보 웹서비스 - 서울특별시 강남구 청담동 도표예보</title>\n'
1606     '<link>http://www.kma.go.kr/weather/main.jsp</link>\n'
1607     '<description>동네예보 웹서비스</description>\n'
1608     '<language>ko</language>\n'
1609     '<generator>동네예보</generator>\n'
1610     '<pubDate>2019년 01월 15일 (화)요일 14:00</pubDate>\n'
1611     '<item>\n'
1612     '<author>기상청</author>\n'
1613     '<category>서울특별시 강남구 청담동</category>\n'
1614     '<title>동네예보(도표) : 서울특별시 강남구 청담동 '
1615     '[X=61,Y=126]</title><link>http://www.kma.go.kr/weather/forecast/timese
1616     ries.jsp?searchType=INTEREST&dongCode=1168056500</link>\n'
1617     '<guid>http://www.kma.go.kr/weather/forecast/timeseries.jsp?searchType=I
1618     NTEREST&dongCode=1168056500</guid>\n'
1619     '<description>\n'
1620     ...
1621     ...
1622
1623     -get method의 params option을 활용하기
1624
1625     api_url = 'http://www.kma.go.kr/wid/queryDFSRSS.jsp'

```



```

1621
1622     payload = {'zone':'1168056500'}
1623
1624     weather_data = requests.get(api_url, payload).text
1625
1626     weather_data
1627     -----
1628     '<?xml version="1.0" encoding="UTF-8" ?>\n<rss
version="2.0">\n<channel>\n<title>기상청 동네예보 웹서비스 - 서울특별시 강남구 청
담동 도표예보
</title>\n<link>http://www.kma.go.kr/weather/main.jsp</link>\n<descriptio
n>동네예보 웹서비스</description>\n<language>ko</language>\n<generator>
동네예보</generator>\n<pubDate>2019년 01월 15일 (화)요일
14:00</pubDate>\n <item>\n<author>기상청</author>\n<category>서울특별
시 강남구 청담동</category>\n<title>동네예보(도표) : 서울특별시 강남구 청담동
[X=61,Y=126]</title><link>http://www.kma.go.kr/weather/forecast/timeseri
es.jsp?searchType=INTEREST&dongCode=1168056500</link>\n<guid>
http://www.kma.go.kr/weather/forecast/timeseries.jsp?searchType=INTERES
T&dongCode=1168056500</guid>\n<description>\n <header>\n
<tm>201901151400</tm>\n <ts>4</ts>\n
1629     ...
1630     ...
1631
1632 5)xml.etree.ElementTree module 사용하기
1633     -Python에서 xml data를 다루기 위한 module
1634
1635     import xml.etree.ElementTree as ET
1636
1637     xml_data = ET.fromstring(weather_data)
1638
1639     for tag in xml_data.iter('data'):
1640         print(tag.find('hour').text + "/" + tag.find('temp').text)
1641     -----
1642     18/-2.0
1643     21/-4.0
1644     24/-5.0
1645     ...
1646     ...
1647
1648     list = []
1649     for tag in xml_data.iter('data'):
1650         dic = {'hour': tag.find('hour').text,
1651             'day' : tag.find('day').text,
1652             'temp' : tag.find('temp').text,
1653             'tmx' : tag.find('tmx').text,
1654             'tmn' : tag.find('tmn').text,
1655             'sky' : tag.find('sky').text,
1656             'pty' : tag.find('pty').text,
1657             'wfKor' : tag.find('wfKor').text,
1658             'wfEn' : tag.find('wfEn').text}
1659     list.append(dic)
1660

```

```

1661     df = pd.DataFrame(list, columns=['hour', 'day', 'temp', 'tmx', 'tmn', 'sky', 'pty',
1662                                'wfKor', 'wfEn'])
1663
1664     df
1665     -----
1666         hour  day temp    tmx    tmn  sky pty wfKor  wfEn
1667     0 18  0   -2.0  -999.0 -999.0  2  0   구름 조금 Partly Cloudy
1668     1 21  0   -4.0  -999.0 -999.0  2  0   구름 조금 Partly Cloudy
1669     2 24  0   -5.0  -999.0 -999.0  1  0   맑음      Clear
1670     3  3  1   -6.0   -1.0  -8.0  1  0   맑음      Clear
1671     4  6  1   -7.0   -1.0  -8.0  1  0   맑음      Clear
1672
1672 6)전자신문 RSS
1673     import pandas as pd
1674     import xml.etree.ElementTree as ET
1675     import requests
1676
1677     api_url = 'http://rss.etnews.com/Section901.xml'
1678
1679     etnews_data = requests.get(api_url).text
1680     etnews_data
1681     -----
1682     '<?xml version="1.0" encoding="utf-8" ?>\n<rss version="2.0">\n\r\n
1683     <channel>\r\n  ...
1684
1685     ...
1686
1687     xml_data = ET.fromstring(etnews_data)
1688     for tag in xml_data.iter('item'):
1689         print(tag.find('title').text + "," + tag.find('pubDate').text)
1690     -----
1691     애플 납품 업체들, 줄줄이 실적 하향... '아이폰 쇼크' 후폭풍,Tue, 15 Jan 2019 17:00:00
1692     +0900
1693     지난해 드론 자격증 취득자 전년비 4배 증가...실효성 지적도,Tue, 15 Jan 2019 17:00:00
1694     +0900
1695     화웨이, 韓 스마트워치 시장 진출,Tue, 15 Jan 2019 17:00:00 +0900
1696     ...
1697     ...
1698
1699     list = []
1700
1701     for tag in xml_data.iter('item'):
1702         dic = {'Title':tag.find('title').text,
1703               'Link':tag.find('link').text,
1704               'Author':tag.find('author').text,
1705               'PubDate':tag.find('pubDate').text,
1706               'Guid':tag.find('guid').text}
1707         list.append(dic)
1708
1709     list
1710     -----
1711     [{'Title': "애플 납품 업체들, 줄줄이 실적 하향... '아이폰 쇼크' 후폭풍",
1712      'Link': 'http://www.etnews.com/20190115000273',

```

```

1709         'Author': '윤건일',
1710         'PubDate': 'Tue, 15 Jan 2019 17:00:00 +0900',
1711         'Guid': '20190115000273'}},
1712         ...
1713         ...
1714
1715     df = pd.DataFrame(list, columns=['Title', 'Link', 'Author', 'PubDate', 'Guid'])
1716     df.head()
1717     -----
1718     ...
1719     ...
1720
1721     df.info()
1722     -----
1723     <class 'pandas.core.frame.DataFrame'>
1724     RangeIndex: 30 entries, 0 to 29
1725     Data columns (total 5 columns):
1726     Title      30 non-null object
1727     Link       30 non-null object
1728     Author     30 non-null object
1729     PubDate    30 non-null object
1730     Guid       30 non-null object
1731     dtypes: object(5)
1732     memory usage: 1.2+ KB
1733
1734
1735 9. Crawling & Scraping
1736 1)Crawling
1737     -Web Site의 data를 그대로 취득하는 것
1738 2)Scraping
1739     -Crawling하여 모든 data에서 필요한 것만 추출하거나 변환하는 처리를 포함.
1740 3)BeautifulSoup4
1741     -Scraping을 위한 module
1742     -일반적으로 HTML tag들이 start tag와 end tag가 서로 pair 되지 않을 경우가 많다.
1743     -pair 되지 않아도 아름답게 처리해 주는 module
1744     -https://www.crummy.com/software/BeautifulSoup/bs4/doc/
1745
1746     from bs4 import BeautifulSoup
1747
1748     html_data = requests.get('https://www.naver.com')
1749
1750     soup = BeautifulSoup(html_data.text, 'html.parser')
1751
1752     soup.title
1753     -----
1754     <title>NAVER</title>
1755
1756 4)Naver 영화 평점 Scraping 하기
1757
1758     from bs4 import BeautifulSoup
1759
1760     html_data =

```

```

requests.get('https://movie.naver.com/movie/point/af/list.nhn?page=1')
1761
1762 soup = BeautifulSoup(html_data.text, 'html.parser')
1763
1764 titles = soup.find_all(class_='movie')
1765
1766 title_list = []
1767 for title in titles:
1768     print(title.text)
1769     -----
1770     대학살의 신
1771     스타워즈: 라스트 제다이
1772     내안의 그놈
1773     말모이
1774     말모이
1775     내안의 그놈
1776     언니
1777     내안의 그놈
1778     존 워
1779     마이 리틀 자이언트
1780
1781 for title in titles:
1782     title_list.append(title.text)
1783
1784 point_list = []
1785 points = soup.find_all(class_='point')
1786 for point in points:
1787     point_list.append(point.text)
1788
1789 review_list = []
1790 reviews = soup.find_all(class_='title')
1791 for review in reviews:
1792     rev = review.text
1793     rev = rev.strip()
1794     rev = rev.replace('\t', '')
1795     rev = rev.replace('\n', '')
1796     rev = rev.replace('신고', '')
1797     review_list.append(rev)
1798
1799 df = pd.DataFrame(title_list, columns=['Title'])
1800 df['Point'] = point_list
1801 df['Review'] = review_list
1802 df
1803     -----
1804     Title                Point    Review
1805 0 대학살의 신              7    대학살의 신자식싸움에 부모등 터진다
1806 1 스타워즈: 라스트 제다이    1    스타워즈: 라스트 제다이어거 보느니 로그원 열번 보는게
    낫다
1807 2 내안의 그놈             10    내안의 그놈재밌어요 유치해도 빵빵터짐 진영 연기 잘하네
    요 라미란과 잘...
1808 3 말모이                 10    말모이후반부에 보고 울었습니다 감동적임
1809 4 말모이                 10    말모이재미를 떠나서 역사는 무조건 10점이다

```

1810	5 내안의 그놈 성형은 다이어트	10	내안의 그놈뽀한스토리이지만 재밌게 잘 보고왔어요최고의
1811	6 언니 개연성도 떨어지고 액션도...	5	언니론다 로우지가 언니 역활했으면 그나마 공감이 됐을듯...
1812	7 내안의 그놈 하계 웃기 좋은 영화	9	내안의 그놈ㅋㅋ재밌었음 생각보다 안정적인 연기력 소소
1813	8 존 워	10	존 워액션의 선두주자 키아누리브스!!
1814	9 마이 리틀 자이언트	7	마이 리틀 자이언트동화보다 더 환상적인 거인

5)Naver 평점 1page부터 100page까지 scraping 하기

```

1818 from bs4 import BeautifulSoup
1819
1820 url = 'https://movie.naver.com/movie/point/af/list.nhn?page='
1821
1822 title_list = []
1823 point_list = []
1824 review_list = []
1825
1826 for pge in range(1, 101):
1827     url = url + str(pge)
1828     print(url)
1829     html_data = requests.get(url)
1830     soup = BeautifulSoup(html_data.text, 'html.parser')
1831     titles = soup.find_all(class_='movie')
1832     points = soup.find_all(class_='point')
1833     reviews = soup.find_all(class_='title')
1834     for title in titles:
1835         title_list.append(title.text)
1836     for point in points:
1837         point_list.append(point.text)
1838     for review in reviews:
1839         rev = review.text
1840         rev = rev.strip()
1841         rev = rev.replace('\t', '')
1842         rev = rev.replace('\n', '')
1843         rev = rev.replace('신고', '')
1844         review_list.append(rev)
1845     url = url.split('=')[0] + '='
1846
1847 df = pd.DataFrame(title_list, columns=['Title'])
1848 df['Point'] = point_list
1849 df['Review'] = review_list
1850
1851 df.info()
1852 -----
1853 <class 'pandas.core.frame.DataFrame'>
1854 RangeIndex: 1020 entries, 0 to 1019
1855 Data columns (total 3 columns):
1856 Title    1020 non-null object
1857 Point    1020 non-null object
1858 Review   1020 non-null object

```

```
1859     dtypes: object(3)
1860     memory usage: 24.0+ KB
```