판다스는 데이터용 분석용 오픈소스 파이썬 라이브러리입니다.

데이터프레임과 시리즈라는 두가지 새로운 자료형을 제공하여, 스프레드 시트 형태의 데이터를 불러와 빠르게조작,정렬,병합 할 수 있습니다. 한마디로 파이썬으로 다루는 엑셀 이라고 보시면 될꺼같습니다.

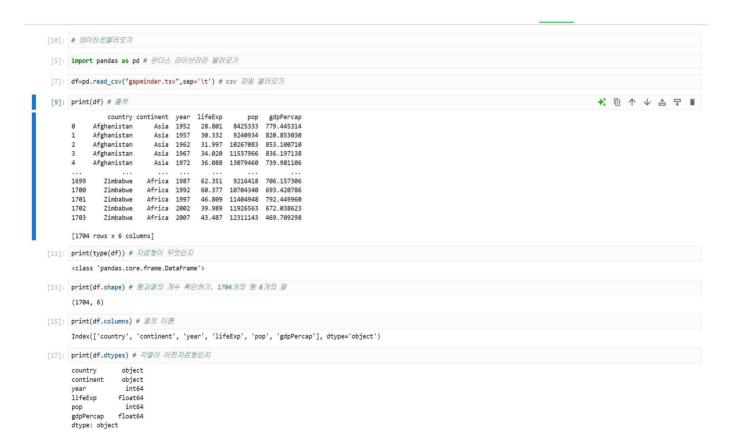
데이터프레임은 전체스프레드시트또는 직사각형의 형태의 데이터로 나타내고,시리즈는 데이터 프레임을 한열을 나타냅니다. 시리즈를 여러개를 모은 딕셔너리나 컬랙션이 판다스의 데이터 프레임이라고 생각해도좋습니다.

데이터를 다룰떼에는 왜 파이썬같은 프로그래밍과 판다스 같은 도구를 사용해야될까요? 첫 번째는 여러 데이터셋에 같은 분석과정을 적용해야 할떼에 일련의 작업을 자동화 할수 있기 때문입니다.

대부분 스프레드시트 프로그램은 마이크로 소프트 365의 VBA와 같은 고유한 메크로 프로그래밍 언어를 제공하지만 이를 활용하는 사람은 많지 않습니다.

메크로보다는 파이썬과 상관없이 작동하는 프로그램언어를 활용하는 것이 좋습니다.

두 번째는 데이터를 작업수행할떼에 데이터에 적용한 모든 실행단계를 기록할수있다는 장점, 즉 재현성이 있기 때문입니다.



```
[]: # 데이터 추출하기
[19]: print(df.head) # 가장 앞 5개의 행을 확인할수있다.
     <bound method NDFrame.head of</pre>
                                      country continent year lifeExp
                                                                        pop gdpPercap
                       Asia 1952 28.801 8425333 779.445314
         Afghanistan
                        Asia 1957 30.332 9240934 820.853030
          Afghanistan
                        Asia 1962 31.997 10267083 853.100710
          Afghanistan
                        Asia 1967 34.020 11537966 836.197138
          Afghanistan
                      Asia 1972 36.088 13079460 739.981106
     4
          Afghanistan
     1699 Zimbabwe Africa 1987 62.351 9216418 706.157306
     1700 Zimbabwe Africa 1992 60.377 10704340 693.420786
     1701
            Zimbabwe Africa 1997 46.809 11404948 792.449960
     1702 Zimbabwe Africa 2002 39.989 11926563 672.038623
     1703 Zimbabwe Africa 2007 43.487 12311143 469.709298
     [1704 rows x 6 columns]>
[21]: country_df=df['country'] # df에서 country 열 데이터를 추출하고 그걸과를country_df 저장
[23]: print(country_df.head()) # country의 가장 앞에 5개의 행 출력
     0 Afghanistan
         Afghanistan
         Afghanistan
         Afghanistan
         Afghanistan
     Name: country, dtype: object
[25]: print(country_df.tail()) # 가장 마지막 5개의 데이터 출력
     1699 Zimbabwe
     1700
           Zimbabwe
     1701
           Zimbabwe
     1702 Zimbabwe
     1703 Zimbabwe
     Name: country, dtype: object
[27]: subset=df[['country','continent','year']] # 3개의 열 데이터를 추출하고 그걸과를 subset변수에 저장
```

```
Afghanistan
Afghanistan
Afghanistan
Afghanistan
Afghanistan
Afghanistan
                                 Zimbabwe
Zimbabwe
Zimbabwe
Zimbabwe
Zimbabwe
                1699
1700
1701
1702
1703
                [1704 rows x 1 columns]
                            country continent year lifeExp pop gdpPercap
Afghanistan Asia 1952 28.801 8425333 779.445314
Afghanistan Asia 1957 30.332 29.40934 820.853030
Afghanistan Asia 1962 31.997 10267083 853.100710
Afghanistan Asia 1967 34.020 11537966 836.197138
Afghanistan Asia 1972 36.088 13070464
 [41]: print(df)
                              Zimbabwe
Zimbabwe
Zimbabwe
Zimbabwe
Zimbabwe
                                                                Africa 1987
Africa 1992
Africa 1997
Africa 2002
Africa 2007
                1699
1700
1701
1702
1703
                                                                                                    62.351 9216418 706.157306
60.377 10704340 693.420786
46.809 11404948 792.449960
39.989 11926563 672.038623
43.487 12311143 469.709298
                [1704 rows x 6 columns]
 [43]: print(df.loc[0]) # 炎世期 谢書与
                country Afghanistan
continent Asia
                                                  Asia
1952
                year
lifeExp
                                                           28.801
                pop 8425333
gdpPercap 779.445314
Name: 0, dtype: object
 [45]: print(df.loc[99]) # 100世期 想書等
                  country
continent
                                              Bangladesh
                                           Bangladesh
Asia
1967
43.453
62821884
721.186086
                  year
lifeExp
                  pop 62821884
gdpPercap 721.186086
Name: 99, dtype: object
 [49]: number_of_row=df.shape[0] # 행계수 구하기
last_row_index=number_of_row=1 # 미지막행 인덱스 구하기
print(df.loc[last_row_index]) # 마지막행의 인덱스로 데이터 추출
                                                   Zimbabwe
Africa
                  country
continent
                   year
lifeExp
                                                    2007
43.487
                   pop 12311143
gdpPercap 469.709298
Name: 1703, dtype: object
  [51]: print(df.tail(n=1)) # 마지막행 추출

        country continent
        year
        lifeExp
        pop
        gdpPercap

        1703
        Zimbabwe
        Africa
        2007
        43.487
        12311143
        469.709298

 [53]: print(df.loc[[0,99,999]]) # 첫번째 100번째 1000번째 데이터 추출

        country continent
        year
        lifeExp
        pop
        gdpPercap

        0
        Afghanistan
        Asia
        1952
        28.881
        8425333
        779.445314

        99
        Bangladesh
        Asia
        1967
        43.435
        62821884
        721.18888

        999
        Mongolia
        Asia
        1967
        51.253
        1149500
        1226.041130

•[35]: print(df.iloc[-1]) # 마지막 행
                                                 Zimbabwe
Africa
2007
43.487
                  country
continent
                   year
lifeExp
                  pop 12311143
gdpPercap 469.709298
Name: 1703, dtype: object
                                                      12311143
```

[39]: print(country_df_list) # 열데이터를 추출 country

```
[37]: subset=df.loc[:,['year','pop']]
print(subset) # :,[曾]] 은 특정 열의 모든 행 출력
                               year pop
1952 8425333
1957 9240934
1962 10267083
1967 11537966
1972 13079460
                 1699 1987 9216418
1700 1992 10704340
1701 1997 11404948
1702 2002 11926563
1703 2007 12311143
                 [1704 rows x 2 columns]
 [39]: subset=df.iloc[:,[2,4,-1]] # 3,5번째와 마지막(-1) 열 데이터의 추출 합니다.
                 print(subset)
                               year pop gdpPercap
1952 8425333 779.445314
1957 9248934 820.853639
1962 10267083 853.160710
1967 11537966 836.197138
1972 13679460 739.981106
                 [1704 rows x 3 columns]
[43]: small_range=list(range(5))
    print(small_range)
                 [0, 1, 2, 3, 4]
             [45]: subset=df.iloc[:,small_range]# 리스트를 사용해서 데이터프레임에서 얼을 추출한
print(subset)
                                          country
Afghanistan
Afghanistan
Afghanistan
Afghanistan
Afghanistan
                                               Zimbabwe
Zimbabwe
Zimbabwe
Zimbabwe
Zimbabwe
                                                                            Africa 1987
Africa 1992
Africa 1997
Africa 2002
Africa 2007
                                                                                                              62.351 9216418
60.377 10704340
46.809 11404948
39.989 11926563
43.487 12311143
                            [1704 rows x 5 columns]
            [47]: small_range=list(range(3,6)) # 3이상 6미만
print(small_range)
                             [3, 4, 5]
             [53]: subset=df.iloc[:,small_range]# 리스트를 사용해서 데이터프레일에서 일을 추출함 print(subset)

        country continent
        year
        lifeExp
        pop

        Afghanistan
        Asia
        1952
        28.801
        8425333

        Afghanistan
        Asia
        1957
        30.332
        9240934

        Afghanistan
        Asia
        1962
        31.997
        10267083

        Afghanistan
        Asia
        1967
        34.020
        11537966

        Afghanistan
        Asia
        1972
        36.088
        13079460

                                                Zimbabwe
Zimbabwe
Zimbabwe
Zimbabwe
Zimbabwe
                                                                            Africa 1987
Africa 1992
Africa 1997
Africa 2002
Africa 2007
                                                                                                              62.351 9216418
60.377 10704340
46.809 11404948
39.989 11926563
43.487 12311143
```

[1704 rows x 5 columns]

```
subset=df.iloc[:,small_range]
                   print(subset)
                          country year pop
Afghanistan 1952 8425333
                         Afghanistan 1957 9240934
                          Afghanistan 1962 10267083
                   3 Afghanistan 1967 11537966
4 Afghanistan 1972 13079460
                   ... 1699 Zimbabwe 1987 9216418
                          Zimbabwe 1987 9216418
Zimbabwe 1992 10704340
Zimbabwe 1997
                   1700
                   1701
                             Zimbabwe 1997 11404948
                             Zimbabwe 2002 11926563
Zimbabwe 2007 12311143
                   1702
                   1703
                   [1704 rows x 3 columns]
         [59]: print(df.columns)
                   Index(['country', 'continent', 'year', 'lifeExp', 'pop', 'gdpPercap'], dtype='object')
         [61]: subset=df.iloc[:, :3]
                   print(subset)
                                country continent year
                        Afghanistan
Afghanistan
                                              Asia 1952
                                                 Asia 1957
                          Afghanistan
                          Afghanistan
                                                Asia 1967
                                              Asia 1972
                          Afghanistan
                   1699
                             Zimbabwe
                                              Africa 1987
                   1700
                              7imbabwe
                                              Africa 1992
                              Zimbabwe
                                              Africa 1997
                   1702
                              Zimbabwe
                                              Africa 2002
                   1703 Zimbabwe Africa 2007
                   [1704 rows x 3 columns]
[63]: subset=df.iloc[:, 0:6:2]

        country
        year
        pop

        Afghanistan
        1952
        8425333

        Afghanistan
        1957
        9240934

        Afghanistan
        1962
        10267083

        Afghanistan
        1967
        11537966

        Afghanistan
        1972
        13079460

                     Zimbabwe 1987 9216418
Zimbabwe 1992 10704340
Zimbabwe 1997 11404948
Zimbabwe 2002 11926563
Zimbabwe 2007 12311143
           1699
           1700
           1701
           1792
           [1704 rows x 3 columns]
[65]: print(df.loc[42,'country']) # country 열에서 행이름이 42인거를 추출 할수있음
          Angola
[67]: print(df.iloc[42,0])
[69]: print(df.loc[[0,99,999],['country','lifeExp','gdpPercap']])
          country lifeExp gdpPercap
0 Afghanistan 28.801 779,445314
99 Bangladesh 43.453 721.186086
999 Mongolia 51.253 1226.041130
```

[57]: small_range=list(range(0,6,2))

```
*[71]: print(df.groupby('year')['lifeExp'].mean()) # 연도별그룹회 LifeExp열 선택 평균계산
                             year
1952
1957
1962
1967
                                            49.057620
51.507401
53.609249
55.678290
57.647386
59.570157
                             1972
                             1977
                            1977 59.576157
1982 61.533197
1987 63.212613
1992 64.168338
1997 65.014676
2002 65.694923
2007 67.007423
Name: lifeExp, dtype: float64
                            multi_group_var=df.groupby(['year','continent'])[['lifeExp','gdpPercap']].mean() # 2 の台 2 コ音劇 print(multi_group_var)
                          year continent
1952 Africa
Americas
Asia
Europe
Oceania
1957 Africa
Americas
Asia
Europe
Oceania
1962 Africa
Americas
Asia
Europe
Oceania
1967 Africa
Americas
Asia
Europe
Oceania
                                                              lifeExp
                                                                                       gdpPercap
                                                           39.135500 1252.572466
53.279840 4079.062552
46.314394 5195.484004
64.408500 5661.057435
69.255000 10298.085650
                                                           69.155600 10298.085550
41.266346 1385.236605
55.960280 4616.043733
49.318544 5787.732940
66.793667 6965.012316
70.259600 11598.522455
58.398760 4901.541870
51.563223 5729.366625
68.539233 8365.468514
71.085600 12696.452430
45.334538 2659.358361
60.410920 5668.253496
54.663640 5971.173374
                                                           60.419926 5668.253496
54.663640 5971.173374
69.737600 10143.823757
71.3180800 14495.021799
74.7459942 2339.615674
62.394920 6491.334139
57.319269 8187.468699
70.775933 12479.575246
71.9180800 16417.333388
                            Asia
Europe
Oceania
1972 Africa
Americas
Asia
Europe
Oceania
•[79]: print(df.groupby('continent')['country'].nunique()) # 그룹화한 데이터개수구하기
                continent
                Africa
Americas
                                         25
33
                Asia
                Europe
                                         30
                Oceania 2
Name: country, dtype: int64
 [83]: print(df.groupby('continent')['country'].value_counts()) # 지정한 열이나 행의개수(빈도수) 구함
               continent country
Africa Algeria
                                       Angola
Benin
                                                                             12
12
                                       Botswana
Burkina Faso
                                                                            12
                                                                             12
                Europe
                                      Switzerland
                                       Turkey
United Kingdom
                                                                             12
                Oceania
                                      Australia
                                                                             12
               New Zealand 12
Name: count, Length: 142, dtype: int64
 [89]: global_yearly_life_expenctancy=df.groupby('year')['lifeExp'].mean() print(global_yearly_life_expenctancy)
                1952
                                49.057620
                               49.057620
51.507401
53.609249
55.678290
57.647386
59.570157
                1957
1962
1967
                1972
                1977
                1982
1987
                                61.533197
63.212613
                1992
                                64.160338
                                65.014676
65.694923
67.007423
                1997
                Name: lifeExp. dtvpe: float64
 [91]: global_yearly_life_expenctancy.plot()
 [91]: <Axes: xlabel='year'>
```

