

Міністерство освіти і науки України Національний технічний університет України "Київський політехнічний інститут імені Ігоря Сікорського" Факультет інформатики та обчислювальної техніки Кафедра автоматики та управління в технічних системах

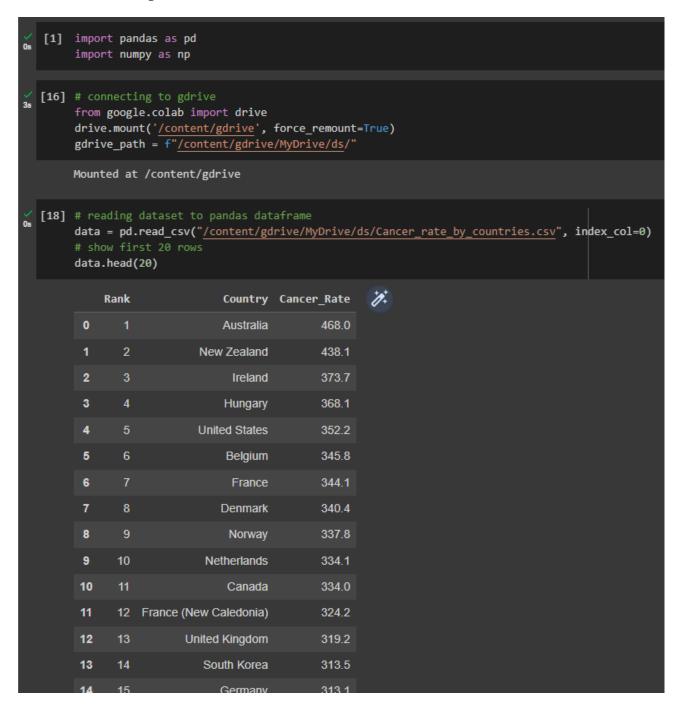
Лабораторна робота №1 Підсумкові функції та відображення

Луцай Катерина	Новотарський М. А.
студентка групи IT-91:	Перевірив:
Виконала	

Мета: навчитися вибрати відповідні дані з DataFrame або Series, вибирати правильні дані з набору даних.

Bapiaнт: 15 – dataset "Cancer rate by countries" (https://www.kaggle.com/dianapratiwi/cancer-rate-by-countries)

Хід виконання роботи:



```
[4] # show dataframe general info
      data.info()
      <class 'pandas.core.frame.DataFrame'>
      Int64Index: 50 entries, 0 to 49
      Data columns (total 3 columns):
       # Column Non-Null Count Dtype
       0 Rank 50 non-null int64
1 Country 50 non-null object
2 Cancer_Rate 50 non-null float64
      dtypes: float64(1), int64(1), object(1)
      memory usage: 1.6+ KB
 Dataset has 3 columns and all rows are filled with values
      data.Cancer Rate
            468.0
     0
 ₽
            438.1
      2
            373.7
            368.1
             352.2
[17] # show columns description
    data.Cancer_Rate.describe()
             50.000000
    count
            294.082000
    mean
            50.026628
            233.600000
    min
    25%
            253.900000
            286.500000
    50%
            317.775000
    75%
    max
            468.000000
    Name: Cancer_Rate, dtype: float64
Column has a numerical data type, thus description shows its size, mean, standard deviation, min and max, and quartiles.
[6] data.Country.describe()
                    50
    unique
                    50
            Australia
    freq
    Name: Country, dtype: object
```

Column has a string data type, thus description shows its size, number of unique values, the most frequent entry, and its

frequency.

[7] # show column mean value
 data.Cancer_Rate.mean()

294.082000000000005

```
[8] # show a list of column's unique values
      data.Country.unique()
      array(['Australia', 'New Zealand', 'Ireland', 'Hungary', 'United States',
               'Belgium', 'France', 'Denmark', 'Norway', 'Netherlands', 'Canada',
               'France (New Caledonia)', 'United Kingdom', 'South Korea',
               'Germany', 'Switzerland', 'Luxembourg', 'Serbia', 'Slovenia', 'Latvia', 'Slovakia', 'Czech Republic', 'Sweden', 'Italy',
               'Croatia', 'Lithuania', 'Estonia', 'Greece', 'Spain', 'Finland', 'Uruguay', 'Belarus', 'Portugal', 'Iceland', 'France (Guadeloupe)',
               'United States (Puerto Rico)', 'Moldova', 'Poland', 'Cyprus', 'France (Martinique)', 'Malta', 'Singapore', 'Japan', 'Austria',
               'Barbados', 'France (French Guiana)', 'Bulgaria', 'Lebanon',
               'France (French Polynesia)', 'Israel'], dtype=object)
[9] # show column's entries and how many times they appear
      data.Country.value counts()
      Australia
      Poland
                                            1
      Greece
                                            1
      Spain
                                            1
      Finland
                                            1
      Uruguay
                                            1
      Relarus
[12] # create a column of mapped values from another column series by substracting
      # column's mean value from all its rows
      cancer mean = data.Cancer Rate.mean()
      data["Cancer div"] = data.Cancer Rate.map(lambda p: p - cancer mean)
      data
       18
              19
                                    Siovenia
                                                      304.9
                                                                   าบ.ชาช
       19
              20
                                       Latvia
                                                      302.2
                                                                    8.118
       20
              21
                                                      297.5
                                                                    3.418
                                    Slovakia
       21
              22
                             Czech Republic
                                                      296.7
                                                                    2.618
       22
             23
                                                      294.7
                                                                    0.618
                                     Sweden
                                                                   -3 482
       23
              24
                                        Italy
                                                      290.6
       24
              25
                                                                   -6.882
                                     Croatia
                                                      287.2
       25
              26
                                    Lithuania
                                                      285.8
                                                                   -8.282
```

Above is Cancer Rate diviation next to Cancer Rate

```
# function to make Country column entries in the upper case, and
    # reset Cancer div to the absolute values (negatives change to positives)
    def remean_points(row):
        row.Cancer_div = np.abs(row.Cancer_Rate - cancer_mean)
        row.Country = row.Country.upper()
        return row
    # applying the defined function to all dataframe rows
    data = data.apply(remean_points, axis=1)
    data
     18
                                                               าบ.ชาช
           19
                                  SLUVENIA
                                                    304.9
₽
     19
           20
                                     LATVIA
                                                    302.2
                                                                8.118
     20
           21
                                                                3.418
                                  SLOVAKIA
                                                    297.5
     21
           22
                           CZECH REPUBLIC
                                                    296.7
                                                                2.618
     22
          23
                                   SWEDEN
                                                    294.7
                                                                0.618
     23
          24
                                       ITALY
                                                    290.6
                                                                3.482
     24
          25
                                                                6.882
                                   CROATIA
                                                    287.2
     25
           26
                                  LITHUANIA
                                                    285.8
                                                                8.282
```

ESTONIA

283.3

10 782

Above is Cancer Rate diviation absolute values next to Cancer Rate [14] # show difference between a Cancer rate series and its mean data.Cancer_Rate - cancer_mean 0 173.918 1 144.018 2 79.618 3 74.018 58.118 4 51.718 50.018 6 46.318

```
[15] # concat Country and stringified Rank values
     data.Country + " " + data.Rank.astype("string")
     0
                                AUSTRALIA 1
     1
                             NEW ZEALAND 2
     2
                                  IRELAND 3
     3
                                 HUNGARY 4
     4
                           UNITED STATES 5
     5
                                  BELGIUM 6
     6
                                   FRANCE 7
                                  DENMARK 8
     8
                                   NORWAY 9
     9
                            NETHERLANDS 10
                                  CANADA 11
     10
     11
                 FRANCE (NEW CALEDONIA) 12
     12
                         UNITED KINGDOM 13
     13
                            SOUTH KOREA 14
     14
                                 GERMANY 15
     15
                            SWITZERLAND 16
     16
                             LUXEMBOURG 17
     17
                                  SERBIA 18
     18
                                SLOVENIA 19
                                  LATVIA 20
     19
     20
                                SLOVAKIA 21
     21
                         CZECH REPUBLIC 22
     22
                                  SWEDEN 23
     23
                                   ITALY 24
     24
                                 CROATIA 25
     25
                              LITHUANIA 26
     26
                                 ESTONIA 27
                                  GREECE 28
     27
     28
                                   SPAIN 29
     29
                                 FINLAND 30
     30
                                 URUGUAY 31
```

Вихідний код у jupyter notebook:

 $\frac{https://colab.research.google.com/drive/13WVuqk6b3vXyhtyhHg5heAE5LD1muZD}{h?usp=sharing}$

Висновки: було розглянуто основні методи бібліотеки pandas на мові Python для вибору відповідних даних із DataFrame або Series, вибору правильних даних із набору даних, такі як: describe(), mean(), unique(), value_counts(), map(), apply() та інші.