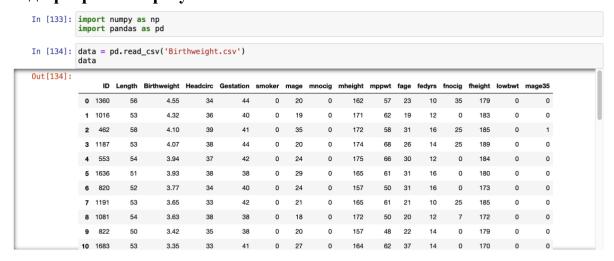
Міністерство освіти і науки України Національний технічний університет України «Київський політехнічний інститут імені Ігоря Сікорського» Факультет інформатики та обчислювальної техніки Кафедра обчислювальної техніки

Лабораторна робота №3

з дисципліни «Аналіз даних з використанням мови Python»

Виконав: студент групи ІП-04 Пащенко Дмитро Олексійович Перевірила: Тимофєєва Ю. С.

Код програми та результат виконання

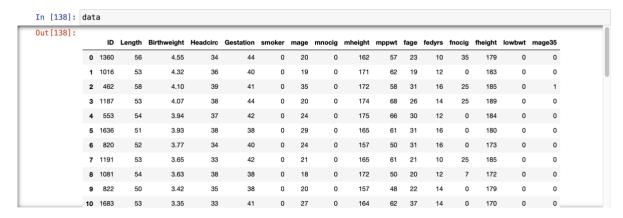


Створити програму, яка за даними файлу відповідно до варіантів лабораторної No2, виконує наступні завдання:

1. Виділити один зі стовпців (на вибір) з файлу як об'єкт Series, виділити з нього підмасив. Задати назви індексів цього об'єкту. Виділити підмасиви за допомогою прямої та непрямої індексацій.

```
In [135]: birthweight = pd.Series(data['Birthweight'], index = [15, 16, 17, 18, 0, 1, 2, 3])
Out[135]:
          15
                3.11
                3.03
          17
                2.92
                2.90
          18
                4.32
                4.10
          Name: Birthweight, dtype: float64
In [136]: birthweight.loc[15:17]
Out[136]: 15
          16
                3.03
                2.92
          Name: Birthweight, dtype: float64
In [137]: birthweight.iloc[0:3]
Out[137]: 15
          16
                3.03
          Name: Birthweight, dtype: float64
```

2. До об'єкту DataFrame, в який записано вміст файлу, додати новий стовпець, що є результатом операцій над іншими стовпцями. Також продемонструвати додавання та видалення рядків, видалення стовпців.



```
In [159]: data_two = data
    data_three = pd.concat([data, data_two], ignore_index = True)
            data_three
                              Birthweight Headcirc smoker mage mnocig mheight mppwt fage
                                                                                                     fnocig
              0 1360
                                    4.55
                                                                             162
                                                                                     57
              2 462
                                                                             172
                                                                                           31
              3
                 1187
                          53
                                    4.07
                                                        0
                                                             20
                                                                      0
                                                                             174
                                                                                     68
                                                                                           26
                                                                                                  14
                                                                                                         25
                                                                                                               189
                 553
                          54
                                    3.94
                                                        0
                                                             24
                                                                      0
                                                                             175
                                                                                     66
                                                                                           30
                                                                                                  12
                                                                                                         0
                                                                                                               184
                                                                                                                                         6
                                    3.41
             77 619
                          52
                                               33
                                                             23
                                                                     25
                                                                             181
                                                                                     69
                                                                                           23
                                                                                                  16
                                                                                                         2
                                                                                                               181
                                                                                                                         0
                                                                                                                                         0
                                                                                                                                 0
                          49
                                                             31
                                                                                     57
                                               34
                                                                     25
                                                                             162
                                                                                           32
                                                                                                  16
                                                                                                               194
             78
                 1369
                                    3.18
                                                                                                         50
                                                                                                                                 0
             79
                 1262
                          53
                                    3.19
                                               34
                                                             27
                                                                     35
                                                                             163
                                                                                     51
                                                                                           31
                                                                                                  16
                                                                                                        25
                                                                                                               185
                                                                                                                                         4
                  516
                          47
                                    2.66
                                               33
                                                             20
                                                                      35
                                                                             170
                                                                                     57
                                                                                           23
                                                                                                  12
                                                                                                         50
                                                                                                               186
                                                                                                                                         3
                          54
                                    4.00
                                                        0
                                                             22
                                                                      0
                                                                             170
                                                                                     53
                                                                                           33
                                                                                                  10
                                                                                                               180
```

3. Встановити один зі стовпців індексом. Визначити основні статистичні характеристики та типи даних всіх стовпців. Змінити тип даних для одного з стовпців. Згрупувати дані за одним зі стовпців, застосувати кілька агрегуючих функцій, виділити підмасив за певними ознаками.

In [160]: data_three = data_three.set_index(['ID']) data_three Out[160]: Length Birthweight Headcirc smoker mage mnocig mheight mppwt fage fedyrs fnocig fheight lowbwt mage35 ageDiff ID 4.55 4.32 4.10 -4 4.07 3.94 3.41 3.18

In [161]: data_three.describe()

Out[161]:

	Length	Birthweight	Headcirc	smoker	mage	mnocig	mheight	mppwt	fage	fedyrs	fnocig	fheight	lowbwt	
count	82.000000	82.000000	82.000000	82.000000	82.000000	82.000000	82.000000	82.000000	82.000000	82.000000	82.000000	82.000000	82.000000	82
mean	51.390244	3.340732	34.585366	0.512195	25.317073	8.439024	164.682927	57.536585	29.121951	13.512195	17.682927	180.707317	0.146341	C
std	2.963798	0.609883	2.413702	0.502927	5.363061	10.809344	6.469068	7.101026	6.839230	2.212445	17.129051	6.914804	0.355623	C
min	43.000000	1.920000	30.000000	0.000000	18.000000	0.000000	149.000000	45.000000	19.000000	10.000000	0.000000	169.000000	0.000000	(
25%	50.000000	3.000000	33.000000	0.000000	21.000000	0.000000	161.000000	53.000000	23.000000	12.000000	0.000000	176.000000	0.000000	(
50%	52.000000	3.320000	34.000000	1.000000	24.000000	2.000000	165.000000	57.000000	30.000000	14.000000	25.000000	181.000000	0.000000	C
75%	53.000000	3.770000	36.000000	1.000000	29.000000	12.000000	170.000000	62.000000	33.000000	16.000000	25.000000	185.000000	0.000000	C
max	58.000000	4.570000	39.000000	1.000000	41.000000	35.000000	181.000000	78.000000	46.000000	16.000000	50.000000	200.000000	1.000000	1

In [162]: data_three.dtypes

Out[162]: Length int64 float64 int64 Birthweight Headcirc smoker int64 int64 mage mnocig int64 int64 mheight mppwt int64 fage fedyrs int64 int64 fnocig int64 fheight int64 lowbwt int64 mage35 int64 ageDiff int64 dtype: object

```
In [165]: data_three.astype({'smoker': 'bool'})
   Out[165]:
                      Length Birthweight Headcirc smoker mage mnocig mheight mppwt fage fedyrs fnocig fheight lowbwt mage35 ageDiff
                  ID
                1360
                         56
                                  4.55
                                             34
                                                  False
                                                          20
                                                                   0
                                                                          162
                                                                                 57
                                                                                      23
                                                                                             10
                                                                                                    35
                                                                                                          179
                                                                                                                   0
                                                                                                                           0
                                                                                                                                  3
                1016
                         53
                                   4.32
                                             36
                                                  False
                                                           19
                                                                   0
                                                                          171
                                                                                  62
                                                                                       19
                                                                                             12
                                                                                                     0
                                                                                                          183
                                                                                                                   n
                                                                                                                           n
                                                                                                                                  0
                 462
                         58
                                   4.10
                                             39
                                                  False
                                                           35
                                                                   0
                                                                          172
                                                                                 58
                                                                                      31
                                                                                             16
                                                                                                    25
                                                                                                          185
                                                                                                                   0
                                                                                                                                  -4
                1187
                         53
                                   4.07
                                             38
                                                  False
                                                           20
                                                                   0
                                                                          174
                                                                                  68
                                                                                      26
                                                                                             14
                                                                                                    25
                                                                                                          189
                                                                                                                   0
                                                                                                                           0
                                                                                                                                  6
                 553
                         54
                                  3.94
                                            37
                                                  False
                                                          24
                                                                   0
                                                                          175
                                                                                 66
                                                                                      30
                                                                                             12
                                                                                                    0
                                                                                                          184
                                                                                                                   0
                         52
                                             33
                                                   True
                                                          23
                                                                  25
                                                                          181
                                                                                 69
                                                                                      23
                                                                                             16
                 619
                                  3.18
                                                          31
                                                                  25
                                                                          162
                                                                                 57
                                                                                      32
                                                                                             16
                                                                                                    50
                1369
                                                   True
                                  3.19
                                            34 True
                                                                         163
                                                                                 51 31
                                                                                             16
                 516 47 2.66 33 True 20 35 170
In [166]: data_three.astype({'smoker': 'bool'}).dtypes
Out[166]: Length
Birthweight
                                int64
                              float64
            Headcirc
                                int64
            smoker
                                  bool
            mage
                                 int64
            mnocig
                                 int64
            mheight
                                 int64
                                 int64
            mppwt
            fage
                                 int64
             fedyrs
                                 int64
             fnocig
                                 int64
                                 int64
             fheight
            lowbwt
                                 int64
            mage35
                                 int64
            ageDiff
dtype: object
                                int64
In [189]: data_three.groupby('smoker').agg([sum, np.min, np.max])
Out[189]:
                                                                                                                            ageDiff
                     Length
                                      Birthweight
                                                        Headcirc
                                                                          mage ... fheight lowbwt
                                                                                                           mage35
                     sum amin amax sum amin amax
             smoker
                                                                           976 ...
                                                                                            2
                                                                                                        1
                                                                                                             2
                                  58 141.54 2.65 4.55 1396
                  0 2080
                            43
                                                               32
                                                                     39
                                                                                      193
                                                                                                  0
                                                                                                                   0
                                                                                                                       1 126 -4
                                                                                                                                         11
                  1 2134
                            46
                                  58 132.40 1.92 4.57 1440
                                                                      39
                                                                                      200
                                                                                                                                          10
                                                                30
                                                                          1100 ...
                                                                                            10
                                                                                                  0
                                                                                                                   0
                                                                                                                          1 166
            2 rows × 42 columns
            4. Створити декілька власних об'єктів DataFrame за такою ж тематикою, що й файл. Наприклад, якщо тема файлу – жаби, можна створити об'єкти, що містять розміри жаб, вагу, стать, кількість особин в популяції і
            т.д. Використати описані в теоретичних відомостях параметри методів merge та concat для різних видів
            злиття та об'єднання даних цих об'єктів.
In [200]: # Один-до-одного
            # OUNN-QO-OUNDO 0
a = pd.DataFrame({'mage': [20, 19, 32, 25], 'Birthweight': [2.4, 4.2, 3.3, 2.7]})
b = pd.DataFrame({'mage': [20, 19, 32, 25], 'smoker': [1, 0, 1, 1]})
c = pd.merge(a, b)
Out[200]:
                mage Birthweight smoker
             0
                  20
                            2.4
```

	mage	Birthweight	smoker
0	20	2.4	1
1	19	4.2	0
2	32	3.3	1
3	25	2.7	1

```
In [201]: # Багато-до-одного
a = pd.DataFrame({'mage': [20, 19, 32, 25], 'smoker': [1, 0, 1, 1], 'Birthweight': [2.4, 4.2, 3.3, 2.7]})
b = pd.DataFrame({'mage': [20, 19, 32, 25], 'fage': [24, 18, 36, 30]})
             c = pd.merge(a, b)
Out[201]:
                mage smoker Birthweight fage
              1 19
             2 32 1 3.3 36
In [202]: # Багато-до-багатьох a = pd.DataFrame({'mage': [20, 19, 32, 25], 'smoker': [1, 0, 1, 1], 'Birthweight': [2.4, 4.2, 3.3, 2.7]}) b = pd.DataFrame({'mage': [20, 19, 32, 25], 'fage': [24, 18, 36, 30], 'Length': [48, 50, 53, 49]})
             c = pd.merge(a, b)
Out[202]:
                mage smoker Birthweight fage Length
              0 20 1 2.4 24 48
              1 19
                            0
                                      4.2 18
                                                    50
             2 32 1 3.3 36 53
              3 25
                        1 2.7 30
                                                    49
 In [203]: a = pd.DataFrame({'mage': [20, 19, 32, 25], 'smoker': [1, 0, 1, 1]})
b = pd.DataFrame({'mage': [21, 19, 33, 26], 'fage': [24, 18, 36, 30]})
c = pd.merge(a, b)
 Out[203]:
                mage smoker fage
              0 19 0 18
In [204]: a = pd.DataFrame({'mage': [20, 19, 32, 25], 'smoker': [1, 0, 1, 1]})
b = pd.DataFrame({'mage': [21, 19, 33, 26], 'fage': [24, 18, 36, 30]})
c = pd.merge(a, b, how = 'outer')
c
 Out[204]:
                 mage smoker fage
              0 20 1.0 NaN
                   19
                           0.0 18.0
              2 32 1.0 NaN
              3 25
                          1.0 NaN
              4 21 NaN 24.0
              5 33 NaN 36.0
              6 26 NaN 30.0
In [205]: c = pd.merge(a, b, how = 'left')
Out[205]:
                mage smoker fage
                   20 1 NaN
             0
                   19
                            0 18.0
              1
             2 32 1 NaN
             3 25 1 NaN
In [206]: c = pd.merge(a, b, how = 'right')
Out[206]:
                mage smoker fage
             0 21 NaN 24
                   19
                          0.0
             2 33 NaN 36
```

```
In [213]: a = pd.DataFrame({'mage': [20, 19, 32, 25], 'smoker': [1, 0, 1, 1]})
    b = pd.DataFrame({'mage': [21, 19, 33, 26], 'smoker': [1, 1, 1, 0]})
    c = pd.concat([a, b])
    c
Out[213]:
                      mage smoker
                  0 20
                  2
                        33 1
                  2
In [214]: a = pd.DataFrame({'mage': [20, 19, 32, 25], 'smoker': [1, 0, 1, 1]})
b = pd.DataFrame({'mage': [21, 19, 33, 26], 'smoker': [1, 1, 1, 0]})
c = pd.concat([a, b], axis = 1)
Out [214]:
                      mage smoker mage smoker
                  0 20 1 21
                          19
                                 0 19
                  1
                  2 32 1 33 1
                  3 25
                                1 26
                                                          0
In [215]: a = pd.DataFrame({'mage': [20, 19, 32, 25], 'smoker': [1, 0, 1, 1]})
b = pd.DataFrame({'mage': [21, 19, 33, 26], 'smoker': [1, 1, 1, 0]})
c = pd.concat([a, b], ignore_index = True)
c
```

Out[215]:

	mage	smoker
0	20	1
1	19	0
2	32	1
3	25	1
4	21	1
5	19	1
6	33	1
7	26	0