

Министерство образования и науки Российской Федерации Федеральное государственное бюджетное образовательное учреждение высшего образования

«Московский государственный технический университет имени Н.Э. Баумана

(национальный исследовательский университет)» (МГТУ им. Н.Э. Баумана)

ФАКУЛЬТЕТ

ИНФОРМАТИКА И СИСТЕМЫ УПРАВЛЕНИЯ

КАФЕДРА СИСТЕМЫ ОБРАБОТКИ ИНФОРМАЦИИ И УПРАВЛЕНИЯ

Отчет по лабораторной работе № 2 «**Изучение библиотек обработки данных**» по курсу "Технологии машинного обучения"

Исполнитель: Студент группы ИУ5-63 Желанкина А.С. 28.02.2019

Задание лабораторной работы

Часть 1.

Выполните первое демонстрационное задание "demo assignment" под названием "Exploratory data analysis with Pandas" со страницы курса https://mlcourse.ai/assignments

Условие задания -

https://nbviewer.jupyter.org/github/Yorko/mlcourse_open/blob/master/jupyter_eng lish/assignments_demo/assignment01_pandas_uci_adult.ipynb?flush_cache=true

Набор данных можно скачать здесь -

https://archive.ics.uci.edu/ml/datasets/Adult

Часть 2.

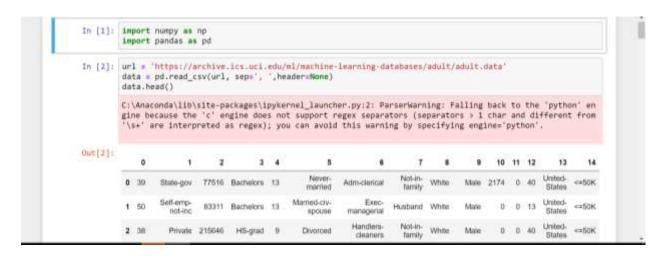
Выполните следующие запросы с использованием двух различных библиотек - Pandas и PandaSQL:

- один произвольный запрос на соединение двух наборов данных
- один произвольный запрос на группировку набора данных с использованием функций агрегирования

Сравните время выполнения каждого запроса в Pandas и PandaSQL.

Экранные формы с текстом программы и примерами её выполнения

ЧАСТЬ 1



```
data.head()
Out[3]:
                                                                                                                             bours.
              age workclass follogt education education-
                                                             marital-
etatus occupation relationship race
                                                                                                             capital-
                                                                                                                     capital-
loss
                                                                                                                               per-
week
                                                                                                               gain
                                                              Never
                                                                           Admi-
           0 39 State-gov 77516 Bachelors
                                                                                                               2174
                                                                                                                                 40
                                                         13
                                                                                 Not-in-family White
                                                                                                       Male
                                                                                                                          0
                                                             married
                                                             Married-
                    Self-emp-
not-inc
                                                                           Exec
               50
                               83311
                                       Bachelors
                                                         13
                                                                                     Husband White
                                                                                                       Male
                                                                                                                  à
                                                                                                                                 13
                                                                      managerial
                                                              spouse
                                                                       Handlers-
              38
                       Private 215646
                                        HS-grad
                                                         9 Divorced
                                                                                 Not-in-family White
                                                                                                       Main
                                                                                                                  0
                                                                                                                          0
                                                                                                                                 40
                                                             Married-
                                                                       Handlers-
            3 53
                       Private 234721
                                                                                     Husband Black
                                                                                                                  0
                                                                                                       Majo
                                                                         cleaners.
                                                              spouse
 In [4]: #How many men and wamen (sex feature) are represented in this dataset?
          print('Total people', data['sex'].count())
print(data['sex'].value_counts())
           Total people 32561
                    21798
10771
           Male
           Female
           Name: sex, dtype: int64
 In [5]: Mihat is the average age (age feature) of women?
           average_age = 0
           counter = 0
           for i in range(data['sex'].count()):
    if data['sex'][i] == 'female':
        counter += 1
                    average_age *= data['age'][1]
           print('The average age of women is ', average_age / counter)
           The average age of women is 36.85823043357163
 In [6]: #What is the percentage of German citizens (native-country feature)?
           counter of german = 0
           for country in data['native-country']:
               if country == 'Germany'
                    counter_of_german ** 1
           print('The percentage of German citizens is '
                  counter_of_german / data['native-country'].count() * 100, '%')
           The percentage of German citizens is 0.42074874850281013 %
In [7]: Midhat are the mean and stundard deviation of age for those who earn more than SBK per year
#(salary feature) and those who earn less than SBK per year?
          ages_of_rich = data.loc[data['salery'] == '550K', 'age']
ages_of_poor = data.loc[data['salery'] == 's550K', 'age']
print("The average age of the rich: (0) +- (1) years.\n\
The average age of the poor: (2) +- (3) years.".format(
    round(ages_of_rich.mean()), round(ages_of_rich.std(), 2),
    round(ages_of_poor.mean()), round(ages_of_poor.std(), 2)))
The average age of the rich: 44 +- 10.52 years.
The average age of the poor: 37 +- 14.02 years.
In [B]: # Is it true that people who earn more than 58K have at least high school education?
           #(education - Bachelors, Prof-school, Assoc-acdm, Assoc-voc, Masters or Doctorate feature)
           data.loc(data['salary'] == '>50K', 'education'].unique()
In [18]: #Display age statistics for each race (race feature) and each gender (sex feature).
           #Use groupby() and describe().
           #Find the maximum age of men of Amer-Indian-Eskima race.
           race_gender_age = data.groupby(['race', 'sex', 'age'])
           print(race_gender_age.describe())
```

```
The maximum age of men of Amer-Indian-Eskimo race is 82
                                             capital-gain
                                                                                      std min
                                                                    mean
                                                     count
                                                                              745.997654 8.8
                                                       2.0 527,500000
           Amer-Indian-Eskimo Female 17
                                        18
                                                              0.000000
                                                                                0.000000 0.0
                                                       2.0
                                                      5.0
                                                                                0.000000
                                        19
                                        21
                                                      4.0
                                                                0.000000
                                                                                0.000000
                                                                                            0.0
                                                      4.0 4307.250000
                                                                             6827.988442
                                        22
                                                                                           0.0
                                                               0.000000
                                                                                0.000000
                                        23
                                                       4.0
                                                                                            0.0
                                        24
                                                      3.0
                                                                0.000000
                                                                                0.000000
                                                                                            0.0
                                        25
                                                      4.0
                                                                0.000000
                                                                                0.000000
                                                                                           0.0
                                        27
                                                       4.0
                                                            831.250000
                                                                             1662.500000
                                                                                            0.0
                                        28
                                                       4.0
                                                                0.000000
                                                                                0.000000
                                        29
                                                       3.0
                                                                0.000000
                                                                                0.000000
                                                                                           0.0
        dick to unumil natput double dick to hide
                                                                0.000000
                                                                                0.000000
                                        30
                                                                                           0.0
                                                       4.0
                                                                0.000000
                                                                                0.000000
                                                                                           0.0
                                        32
                                                       1.0
                                                                0.000000
                                                                                      MaN
                                                                                           0.0
                                                                                0.000000
                                                                0.000000
                                        33
                                                       2.0
                                                                                           0.0
In [11]: #Among whom is the proportion of those who earn a Lat (>50K) greater; married or single men
           #(marital-status feature)? Consider as married those who have a marital-status starting #with Married (Married-civ-spouse, Married-spouse-absent or Married-AF-spouse),
           #the rest are considered bachelors.
          print('Married men:\n')
print(data.loc[(data['sex'] == 'Male') & (data['marital-status'].str.startswith('Married')),
                     'salary'].value_counts())
          Married men:
          <=50%
                   7576
           >58K
                     5965
           Name: salary, dtype: int64
          Single men:
                  7552
           4±5800
           >58K
                      697
          Name: salary, dtype: int64
In [16]: #What is the maximum number of hours a person works per week (hours-per-week feature)? 
Whow many people work such a number of hours and what is the percentage of those who earn
           do Lot among them?
           maximum = data['hours-per-week'].max()
           print('The maximum number of hours a person works per week is', maximum)
           hard = data[data['hours-per-week'] == maximum].shape[0]
           print(hard, 'people work such a number of hours')
           rich = float(data[(data['hours-per-week'] == maximum) &
          (data[ salary'] == >50%')).shape[0]) / hard
print('The percentage of those who earn a lot among them is', rich * 100, '%')
           The maximum number of hours a person works per week is 99
           BS people work such a number of hours
           The percentage of those who earn a lot among them is 29.411764705882355 %
In [26]: #Count the average time of work (hours-per-week) for those who earn a little and a Lot (salary)
           #for each country (native-country). What will these be for Japan!
          little = []
           lot = []
           country_salary = data.groupby(['native-country', 'salary'])
for (country, salary), hours in country_salary:
    print(country, salary, int(hours['hours-per-week'].mean()))
           ? <=50K 40
           ? >50K 45
          Cambodia <=50K 41
Cambodia >50K 40
          Canada <=50K 37
Canada >50K 45
           China <=50K 37
           China >50K 38
          Columbia <=50K 38
Columbia >50K 50
           Cuba <=50K 37
           Cuba >50K 42
          Dominican-Republic <=50K 42
           Dominican-Republic >50K 47
```

```
[ ] [pip3 install pandasql
         [ Collecting pandasql
                 Downloading https://files.pythonhosted.org/packages/6b/c4/ee4096ffa2eeeca0c749b26f8371bd26aa5c8b611c43de99a4f86d
               Requirement already satisfied: numpy in /usr/local/lib/python3.6/dist-packages (from pandasql) (1.14.6)
               Requirement already satisfied: pandas in /usr/local/lib/python3.5/dist-packages (from pandasql) (8.22.8)
               Requirement already satisfied: sqlalchemy in /usr/local/lib/python3.6/dist-packages (from pandasql) (1.2.18)
               Requirement already satisfied: python-dateutil>=2 in /usr/local/lib/python3.6/dist-packages (from pandas->pandasql
               Requirement already satisfied: pytz>=2011k in /usr/local/lib/python3.6/dist-packages (from pandas->pandasql) (2018
               Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.6/dist-packages (from python-dateutil>=2->pandas
              Building wheels for collected packages: pandasql
                 Building wheel for pandasql (setup.py) ... done
                 Stored in directory: /root/.cache/pip/wheels/53/6c/18/b87a2e5fa8a82e9c826311de56218b8d1c81846e18a9687fc9
               Successfully built pandasql
               Installing collected packages: pandasql
              Successfully installed pandasql-0.7.3
               import numpy as np
import pandas as po
import pandasql as
>
        [28]
               from time import time
               url = "https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/wdbc.data"
        [29]
               data_diagnostic = pd.read_csv(url, sep=',',header=None)
data_diagnostic.columns = ['id'] + ['D' + str(i) for i in range(1, 32)]
data_diagnostic.head()
         D+
                          id D1
                                             03
                                                      Da
                                                               DS.
                                                                         06
                                                                                   07
                                                                                                                 022
                                                                                                                                          D25
                                                                                                                                                   D:
                                                                                            DB
                                                                                                                         D23
                                                                                                                                  D24
                     842302 M 17.99 10.38 122.80 1001.0 0.11840 0.27760 0.3001 0.14710
                                                                                                           ... 25.38 17.33 184.60 2019.0 0.16;
                     842517 M 20.57 17.77 132.90 1326.0 0.08474 0.07864 0.0869 0.07017
                                                                                                            ... 24.99 23.41
                                                                                                                              158.80
               2 84300903 M 19.69 21.25 130.00 1203.0 0.10960 0.15990 0.1974 0.12790
                                                                                                            ... 23.57 25.53 152.50
                                                                                                                                       1709.0 0.14
                3 84348301 M 11.42 20.38
                                                  77.58 386.1 0.14250 0.28390 0.2414 0.10520
                                                                                                           ... 14,91 26.50
                                                                                                                               98.87
                                                                                                                                        567.7 0.201
               4 84358402 M 20.29 14.34 135.10 1297.0 0.10030 0.13280 0.1980 0.10430 ... 22.54 16.67 152.20 1575.0 0.13:
               5 rows * 32 columns
              url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/wpbc.data'
data_prognostic = pd.read_csv(url, sep=',',header=Wone)
data_prognostic.columns = ['id'] + ['P' + str(i) for i in range(1, 35)]
data_prognostic.head()
        [30]
         E+
                        id Pl
                                         P3
                                                 P4
                                                         PS.
                                                                  96
                                                                            P7
                                                                                     PB
                                                                                             P9 ...
                                                                                                           P25
                                                                                                                   P26
                                                                                                                            P27
                                                                                                                                     P28
                                                                                                                                              P29
                                 31 18.02 27.60 117.50 1013.0 0.09489 0.1036 0.1086
                                                                                                  ... 139.70 1436.0 0.1195 0.1926 0.3140 0
                     8423 N
                                  61 17.99 10.38 122.80
                                                             1001.0 0.11840 0.2776 0.3001
                                                                                                        184.60 2019.0 0.1622 0.6656
                                116 21.37 17.44 137.50
                                                             1373.0 0.08836 0.1189 0.1255
                                                                                                        159.10
                                                                                                               1949.0 0.1188 0.3449 0.3414 0
               3 843483 N 123 11.42 20.38
                                                      77.58
                                                              386.1 0.14250 D.2839 0.2414
                                                                                                         98.87
                                                                                                                 567.7 0.2098 0.8663 0.6869 0
               4 843584 R 27 20.29 14.34 135.10 1297.0 0.10030 0.1328 0.1980
                                                                                                  ... 152.20 1575.0 0.1374 0.2050 0.4000 0
              5 rows = 35 columns
        [31] # I - diagnosis, 2 - radius, 22 - worst radius
to diagnosis = time()
aggregations_diagnosis = {
    'D22': lambda *: max(x)
>
               print(data_diagnostic.groupby('D1').agg(aggregations_diagnostic))
t1_diagnostic + time()
print('It takes) ', t1_diagnostic - t0_diagnostic)
                     D22
         D:
              D1
                   19.82
                   36.04
              It takes: 0.005372047424316406
               # 1 - diagnosis, 2 - radius, 22 - worst radius

## diagnostic = time()

aggr_query_diagnostic = ""

SLECT
                       D1,
#MH(D22)
                   FROM data diagnostic
GROUP BY Di
```

```
print(ps.sqldf(aggr_query_diagnostic, locals()))
ti_diagnostic ~ time()
print('It takess ', ti_diagnostic - t0_diagnostic)
          0
                D1 max(D22)
          D.
               0 B 19.82
1 M 36.04
               It (diagnostic) takes: 0.04991793632507324
        t1 = time()
print(result)
print('It takes: ', t1 - t0)
                          1d D1
                                     D2 D22 P1 P2
          D.
                     842517 M 28.57 24.99 N 116
843786 M 12.45 15.47 R 77
               0
                      844359 M 18.25 22.88 N 60
844981 M 13.00 15.49 N 110
845636 M 16.02 19.19 N 123
        [36] t0 = time()
merge query = ...
SELECT
d.10,
d.01,
d.02,
d.022,
>
                         p.id,
p.P1,
p.P2
                    FROM data diagnostic as d 30IN data_prognostic as p ON (d.id = p.id) dROUP BV p.id
                print(ps.sqldf(merge_query, locals()))
               t1 = time()
print('It takes: ', t1 - t0)
                                                       id P1 P2
B5715 N 111
                       id D1 D2 D22
85715 M 13.17 15.67
                                       D2 D22
          D.
               a
                        86208 M 20.26 24.22
                                                        86288 R
                                                                     10
                        86517 M 18.66 22.25
                                                        86517 N 108
                        87163 M 13.43 17.98
                                                        87163 N 84
                        87154 M 15.46 18.79
                                                        87164 N 98
                        87880 M 13.81 19.20
89122 M 19.40 23.79
                                                        87888 N 38
                                                        89122 N
                                                                      17
```