Московский государственный технический университет им. Н.Э. Баумана

Факультет «Информатика и системы управления»

Кафедра «Системы обработки информации и управления»

**Отчёт**

**"Методы машинного обучения"**

**Лабораторная работа № 2**

**"Изучение библиотек обработки данных"**

ИСПОЛНИТЕЛЬ:

Студент группы ИУ5-21М

Коростелёв В. М. _____

ПРЕПОДАВАТЕЛЬ:

Гапанюк Ю. Е. _____

**Москва – 2019**

# Задание

Часть 1. Выполните первое демонстрационное задание "demo assignment" под названием "Exploratory data analysis with Pandas" со страницы курса https://mlcourse.ai/assignments (https://mlcourse.ai/assignments) 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.

```
In [1]: import numpy as np
        import pandas as pd
```

```
In [2]: data = pd.read_csv('adult.data.csv')
        data.head()
```

Out[2]:

| | age | workclass | fnlwgt | education | education-num | marital-status | occupation | relationship | race | sex | cap |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 39 | State-gov | 77516 | Bachelors | 13 | Never-married | Adm-clerical | Not-in-family | White | Male | 2 |
| 1 | 50 | Self-emp-not-inc | 83311 | Bachelors | 13 | Married-civ-spouse | Exec-managerial | Husband | White | Male | |
| 2 | 38 | Private | 215646 | HS-grad | 9 | Divorced | Handlers-cleaners | Not-in-family | White | Male | |
| 3 | 53 | Private | 234721 | 11th | 7 | Married-civ-spouse | Handlers-cleaners | Husband | Black | Male | |
| 4 | 28 | Private | 338409 | Bachelors | 13 | Married-civ-spouse | Prof-specialty | Wife | Black | Female | |

```
In [3]: data['sex'].value_counts()
```

```
Out[3]: Male      21790
        Female    10771
        Name: sex, dtype: int64
```

```
In [4]: data.loc[data['sex'] == 'Female', 'age'].mean()
```

```
Out[4]: 36.85823043357163
```

```
In [5]:  print("{0:%}".format(data[data["native-country"] == "Germany"] .shape[0] / data.shape[0]))
```

```
0.420749%
```

```
In [6]:  ages1 = data[data["salary"] == "<=50K"]["age"]
         ages2 = data[data["salary"] == ">50K"]["age"]
         print("<=50K: = {0} ± {1} years".format(ages1.mean(), ages1.std()))
         print(" >50K: = {0} ± {1} years".format(ages2.mean(), ages2.std()))

         <=50K: = 36.78373786407767 ± 14.02008849082488 years
          >50K: = 44.24984058155847 ± 10.519027719851826 years
```

```
In [7]:  high_educations = set(["Bachelors", "Prof-school", "Assoc-acdm", "Assoc-voc", "Master
         s", "Doctorate"])
         def high_educated(e):
             return e in high_educations
         data[data["salary"] == ">50K"]["education"].map(high_educated).all()
```

Out[7]:  False

```
In [8]:  data.groupby(["race", "sex"])["age"].describe()
```

Out[8]:

| | | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|---|
| **race** | **sex** | | | | | | | | |
| **Amer-Indian-Eskimo** | **Female** | 119.0 | 37.117647 | 13.114991 | 17.0 | 27.0 | 36.0 | 46.00 | 80.0 |
| | **Male** | 192.0 | 37.208333 | 12.049563 | 17.0 | 28.0 | 35.0 | 45.00 | 82.0 |
| **Asian-Pac-Islander** | **Female** | 346.0 | 35.089595 | 12.300845 | 17.0 | 25.0 | 33.0 | 43.75 | 75.0 |
| | **Male** | 693.0 | 39.073593 | 12.883944 | 18.0 | 29.0 | 37.0 | 46.00 | 90.0 |
| **Black** | **Female** | 1555.0 | 37.854019 | 12.637197 | 17.0 | 28.0 | 37.0 | 46.00 | 90.0 |
| | **Male** | 1569.0 | 37.682600 | 12.882612 | 17.0 | 27.0 | 36.0 | 46.00 | 90.0 |
| **Other** | **Female** | 109.0 | 31.678899 | 11.631599 | 17.0 | 23.0 | 29.0 | 39.00 | 74.0 |
| | **Male** | 162.0 | 34.654321 | 11.355531 | 17.0 | 26.0 | 32.0 | 42.00 | 77.0 |
| **White** | **Female** | 8642.0 | 36.811618 | 14.329093 | 17.0 | 25.0 | 35.0 | 46.00 | 90.0 |
| | **Male** | 19174.0 | 39.652498 | 13.436029 | 17.0 | 29.0 | 38.0 | 49.00 | 90.0 |

```
In [9]:  data[(data["race"] == "Amer-Indian-Eskimo")
              & (data["sex"] == "Male")]["age"].max()
```

Out[9]:  82

```
In [10]: def is_married(m):
             return m.startswith("Married")
         data["married"] = data["marital-status"].map(is_married)
         (data[(data["sex"] == "Male") & (data["salary"] == ">50K")]
             ["married"].value_counts())
```

Out[10]: True     5965
         False     697
         Name: married, dtype: int64

```
In [11]:  m = data["hours-per-week"].max()
          print("Maximum is {} hours/week.".format(m))
          people = data[data["hours-per-week"] == m]
          c = people.shape[0]
          print("{} people work this time at week.".format(c))
          s = people[people["salary"] == ">50K"].shape[0]
          print("{0:%} get >50K salary.".format(s / c))
```

```
Maximum is 99 hours/week.
85 people work this time at week.
29.411765% get >50K salary.
```

```
In [12]:  p = pd.crosstab(data["native-country"], data["salary"],
                          values=data['hours-per-week'], aggfunc="mean")
          p
```

Out[12]:

| salary | <=50K | >50K |
|---|---|---|
| native-country | | |
| ? | 40.164760 | 45.547945 |
| Cambodia | 41.416667 | 40.000000 |
| Canada | 37.914634 | 45.641026 |
| China | 37.381818 | 38.900000 |
| Columbia | 38.684211 | 50.000000 |
| Cuba | 37.985714 | 42.440000 |
| Dominican-Republic | 42.338235 | 47.000000 |
| Ecuador | 38.041667 | 48.750000 |
| El-Salvador | 36.030928 | 45.000000 |
| England | 40.483333 | 44.533333 |
| France | 41.058824 | 50.750000 |
| Germany | 39.139785 | 44.977273 |
| Greece | 41.809524 | 50.625000 |
| Guatemala | 39.360656 | 36.666667 |
| Haiti | 36.325000 | 42.750000 |
| Holand-Netherlands | 40.000000 | NaN |
| Honduras | 34.333333 | 60.000000 |
| Hong | 39.142857 | 45.000000 |
| Hungary | 31.300000 | 50.000000 |
| India | 38.233333 | 46.475000 |
| Iran | 41.440000 | 47.500000 |
| Ireland | 40.947368 | 48.000000 |
| Italy | 39.625000 | 45.400000 |
| Jamaica | 38.239437 | 41.100000 |
| Japan | 41.000000 | 47.958333 |
| Laos | 40.375000 | 40.000000 |
| Mexico | 40.003279 | 46.575758 |
| Nicaragua | 36.093750 | 37.500000 |
| Outlying-US(Guam-USVI-etc) | 41.857143 | NaN |
| Peru | 35.068966 | 40.000000 |
| Philippines | 38.065693 | 43.032787 |
| Poland | 38.166667 | 39.000000 |
| Portugal | 41.939394 | 41.500000 |
| Puerto-Rico | 38.470588 | 39.416667 |
| Scotland | 39.444444 | 46.666667 |
| South | 40.156250 | 51.437500 |
| Taiwan | 33.774194 | 46.800000 |
| Thailand | 42.866667 | 58.333333 |
| Trinadad&Tobago | 37.058824 | 40.000000 |
| United-States | 38.799127 | 45.505369 |

| | salary | <=50K | >50K |
|---|---|---|---|
| native-country | | | |
| Vietnam | | 37.193548 | 39.200000 |
| Yugoslavia | | 41.600000 | 49.500000 |

In [13]:
```python
p.loc["Japan"]
```

Out[13]:
```
salary
<=50K    41.000000
>50K     47.958333
Name: Japan, dtype: float64
```

In [14]:
```python
from pandasql import sqldf
pysqldf = lambda q: sqldf(q, globals())
```

In [15]:
```python
wind = (pd.read_csv('wind speed.csv', header=None,
                    names=["row", "UNIX", "date",
                           "time", "speed", "text"])
        .drop("text", axis=1))
temp = (pd.read_csv('temperature.csv', header=None,
                    names=["row", "UNIX", "date",
                           "time", "temperature", "text"])
        .drop("text", axis=1))
```

In [16]:
```python
wind.head()
```

Out[16]:

| | row | UNIX | date | time | speed |
|---|---|---|---|---|---|
| 0 | 1 | 1475315718 | 2016-09-30 | 23:55:18 | 7.87 |
| 1 | 2 | 1475315423 | 2016-09-30 | 23:50:23 | 7.87 |
| 2 | 3 | 1475315124 | 2016-09-30 | 23:45:24 | 9.00 |
| 3 | 4 | 1475314821 | 2016-09-30 | 23:40:21 | 13.50 |
| 4 | 5 | 1475314522 | 2016-09-30 | 23:35:22 | 15.75 |

In [17]:
```python
wind.dtypes
```

Out[17]:
```
row        int64
UNIX       int64
date      object
time      object
speed    float64
dtype: object
```

In [18]:
```python
temp.head()
```

Out[18]:

| | row | UNIX | date | time | temperature |
|---|---|---|---|---|---|
| 0 | 1 | 1475315718 | 2016-09-30 | 23:55:18 | 48 |
| 1 | 2 | 1475315423 | 2016-09-30 | 23:50:23 | 48 |
| 2 | 3 | 1475315124 | 2016-09-30 | 23:45:24 | 48 |
| 3 | 4 | 1475314821 | 2016-09-30 | 23:40:21 | 48 |
| 4 | 5 | 1475314522 | 2016-09-30 | 23:35:22 | 48 |

```
In [19]:  temp.dtypes
```

```
Out[19]:  row             int64
          UNIX            int64
          date           object
          time           object
          temperature     int64
          dtype: object
```

```
In [20]:  wind.merge(temp[["UNIX", "temperature"]], on="UNIX").head()
```

Out[20]:

|   | row | UNIX | date | time | speed | temperature |
|---|-----|------|------|------|-------|-------------|
| 0 | 1 | 1475315718 | 2016-09-30 | 23:55:18 | 7.87 | 48 |
| 1 | 2 | 1475315423 | 2016-09-30 | 23:50:23 | 7.87 | 48 |
| 2 | 3 | 1475315124 | 2016-09-30 | 23:45:24 | 9.00 | 48 |
| 3 | 4 | 1475314821 | 2016-09-30 | 23:40:21 | 13.50 | 48 |
| 4 | 5 | 1475314522 | 2016-09-30 | 23:35:22 | 15.75 | 48 |

```
In [21]:  %%timeit
          wind.merge(temp[["UNIX", "temperature"]], on="UNIX")
```

```
          20.6 ms ± 2.41 ms per loop (mean ± std. dev. of 7 runs, 10 loops each)
```

```
In [22]:  pysqldf("""SELECT w.row, w.UNIX, w.date, w.time,
                  w.speed, t.temperature
                  FROM wind AS w JOIN temp AS t
                  ON w.UNIX = t.UNIX """).head()
```

Out[22]:

|   | row | UNIX | date | time | speed | temperature |
|---|-----|------|------|------|-------|-------------|
| 0 | 1 | 1475315718 | 2016-09-30 | 23:55:18 | 7.87 | 48 |
| 1 | 2 | 1475315423 | 2016-09-30 | 23:50:23 | 7.87 | 48 |
| 2 | 3 | 1475315124 | 2016-09-30 | 23:45:24 | 9.00 | 48 |
| 3 | 4 | 1475314821 | 2016-09-30 | 23:40:21 | 13.50 | 48 |
| 4 | 5 | 1475314522 | 2016-09-30 | 23:35:22 | 15.75 | 48 |

```
In [23]:  %%timeit
          pysqldf("""SELECT w.row, w.UNIX, w.date, w.time,
                  w.speed, t.temperature
                  FROM wind AS w JOIN temp AS t
                  ON w.UNIX = t.UNIX """).head()
```

```
          758 ms ± 17.6 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
```

```
In [24]:  wind.groupby("date")["speed"].mean().head()
```

```
Out[24]:  date
          2016-09-01    6.396560
          2016-09-02    5.804086
          2016-09-03    4.960248
          2016-09-04    5.184571
          2016-09-05    5.830676
          Name: speed, dtype: float64
```

```
In [25]: %%timeit
         wind.groupby("date")["speed"].mean().head()
```

3.36 ms ± 75.6 μs per loop (mean ± std. dev. of 7 runs, 100 loops each)

```
In [26]: pysqldf("""SELECT date, AVG(speed) FROM wind GROUP BY date """).head()
```

Out[26]:

|   | date | AVG(speed) |
|---|------|-----------|
| 0 | 2016-09-01 | 6.396560 |
| 1 | 2016-09-02 | 5.804086 |
| 2 | 2016-09-03 | 4.960248 |
| 3 | 2016-09-04 | 5.184571 |
| 4 | 2016-09-05 | 5.830676 |

```
In [27]: %%timeit
         pysqldf("""SELECT date, AVG(speed) FROM wind GROUP BY date """).head()
```

293 ms ± 8.32 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)