

Лабораторная работа (Изучение библиотек обработки данных.)

Часть 1

Создание окружения

```
1 import tensorflow as tf
2 tf.test.gpu_device_name()
```



```
1 !pip install -U -q PyDrive
2
3 from pydrive.auth import GoogleAuth
4 from pydrive.drive import GoogleDrive
5 from google.colab import auth
6 from oauth2client.client import GoogleCredentials
7
8 # 1. Authenticate and create the PyDrive client.
9 auth.authenticate_user()
10 gauth = GoogleAuth()
11 gauth.credentials = GoogleCredentials.get_application_default()
12 drive = GoogleDrive(gauth)
```

```
1 !pip install -U pandasql
```



This file was updated remotely or in another tab. To force a save, overwriting the last update, select Save from the File menu



```
1 file_list = drive.ListFile({'q': "'1FV3XZb9X25mBGBhYOsmg5eyeuJKSnNLp' in parer
2 for file1 in file_list:
3     print('title: %s, id: %s' % (file1['title'], file1['id']))
```



```
1 train_downloaded = drive.CreateFile({'id': '1Se6oWEujEmZtmnJvcDu9NGKXfG7_DX-v'})
2 train_downloaded.GetContentFile('winter.csv')
3 test_downloaded = drive.CreateFile({'id': '117EvqN7ZQtIfntEWtdz42OjFnoNdns1R'})
4 test_downloaded.GetContentFile('summer.csv')
5 train_downloaded = drive.CreateFile({'id': '1CEYqFeTOOPPVIYuZZ9otX7_FWP0FE-e'})
6 train_downloaded.GetContentFile('dictionary.csv')
7 test_downloaded = drive.CreateFile({'id': '11N8YVg6fMtyhlH2D8M796Tkq8y--BmS7'})
8 test_downloaded.GetContentFile('adult.data.csv')
```

```
1 !pip install -q matplotlib-venn
```

```
1 !apt-get -qq install -y libfluidsynth1
```



```
1 # To determine which version you're using:
2 !pip show tensorflow
3
4 # For the current version:
5 !pip install --upgrade tensorflow
6
7 # For a specific version:
8 !pip install tensorflow==1.2
9
10 # For the latest nightly build:
11 !pip install tf-nightly
```



This file was updated remotely or in another tab. To force a save, overwriting the last update, select Save from the File menu



This file was updated remotely or in another tab. To force a save, overwriting the last update, select
Save from the File menu



```
1 # https://pypi.python.org/pypi/libarchive
```

This file was updated remotely or in another tab. To force a save, overwriting the last update, select Save from the File menu

```
↳ Selecting previously unselected package libarchive-dev:amd64.
   (Reading database ... 131309 files and directories currently installed.)
   Preparing to unpack .../libarchive-dev_3.2.2-3.1ubuntu0.3_amd64.deb ...
   Unpacking libarchive-dev:amd64 (3.2.2-3.1ubuntu0.3) ...
   Processing triggers for man-db (2.8.3-2ubuntu0.1) ...
   Setting up libarchive-dev:amd64 (3.2.2-3.1ubuntu0.3) ...
   Building wheel for libarchive (setup.py) ... done
```

```
1 # https://pypi.python.org/pypi/pydot
2 !apt-get -qq install -y graphviz && pip install -q pydot
3 import pydot
```

```
1 !apt-get -qq install python-cartopy python3-cartopy
```

```
2 | import cartopy
```

```
☞ Selecting previously unselected package python-pkg-resources.  
(Reading database ... 131365 files and directories currently installed.)  
Preparing to unpack .../0-python-pkg-resources_39.0.1-2_all.deb ...  
Unpacking python-pkg-resources (39.0.1-2) ...  
Selecting previously unselected package python-pyshp.  
Preparing to unpack .../1-python-pyshp_1.2.12+ds-1_all.deb ...  
Unpacking python-pyshp (1.2.12+ds-1) ...  
Selecting previously unselected package python-shapely.  
Preparing to unpack .../2-python-shapely_1.6.4-1_amd64.deb ...  
Unpacking python-shapely (1.6.4-1) ...  
Selecting previously unselected package python-six.  
Preparing to unpack .../3-python-six_1.11.0-2_all.deb ...  
Unpacking python-six (1.11.0-2) ...  
Selecting previously unselected package python-cartopy:amd64.  
Preparing to unpack .../4-python-cartopy_0.14.2+dfsg1-2build3_amd64.deb ...  
Unpacking python-cartopy:amd64 (0.14.2+dfsg1-2build3) ...  
Selecting previously unselected package python3-pkg-resources.  
Preparing to unpack .../5-python3-pkg-resources_39.0.1-2_all.deb ...  
Unpacking python3-pkg-resources (39.0.1-2) ...  
Selecting previously unselected package python3-pyshp.  
Preparing to unpack .../6-python3-pyshp_1.2.12+ds-1_all.deb ...  
Unpacking python3-pyshp (1.2.12+ds-1) ...  
Selecting previously unselected package python3-shapely.  
Preparing to unpack .../7-python3-shapely_1.6.4-1_amd64.deb ...  
Unpacking python3-shapely (1.6.4-1) ...  
Selecting previously unselected package python3-six.  
Preparing to unpack .../8-python3-six_1.11.0-2_all.deb ...  
Unpacking python3-six (1.11.0-2) ...  
Selecting previously unselected package python3-cartopy:amd64.  
Preparing to unpack .../9-python3-cartopy_0.14.2+dfsg1-2build3_amd64.deb ...  
Unpacking python3-cartopy:amd64 (0.14.2+dfsg1-2build3) ...  
Setting up python-shapely (1.6.4-1) ...  
Setting up python-pyshp (1.2.12+ds-1) ...  
Setting up python3-six (1.11.0-2) ...  
Setting up python3-shapely (1.6.4-1) ...  
Setting up python3-pyshp (1.2.12+ds-1) ...  
Setting up python3-pkg-resources (39.0.1-2) ...  
Setting up python-pkg-resources (39.0.1-2) ...  
Setting up python-six (1.11.0-2) ...  
Setting up python3-cartopy:amd64 (0.14.2+dfsg1-2build3) ...  
Setting up python-cartopy:amd64 (0.14.2+dfsg1-2build3) ...
```

This file was updated remotely or in another tab. To force a save, overwriting the last update, select Save from the File menu

```
1 | import numpy as np  
2 | import pandas as pd
```

```
1 | data = pd.read_csv('adult.data.csv')
```

Работа с данными

```
1 | data.head()
```



1. Сколько мужчин и женщин представлены в этом наборе данных?

```
1 data['sex'].value_counts()
```



2. Каков средний возраст женщин?

```
1 data.loc[data['sex'] == 'Female', 'age'].mean()
```



3. Какова доля немецких граждан?

```
1 float((data['native-country'] == 'Germany').sum()) / data.shape[0]
```



4-5. Каковы среднее значение и стандартное отклонение возраста тех, кто получает более 50К в год (функция зарплаты) и тех, кто получает менее 50К в год?

This file was updated remotely or in another tab. To force a save, overwriting the last update, select Save from the File menu

```
3 print("The average age of the rich: {0} +- {1} years, poor - {2} +- {3} years.  
4     round(ages1.mean()), round(ages1.std(), 1),  
5     round(ages2.mean()), round(ages2.std(), 1))
```



6. Правда ли, что люди, получающие более 50 тысяч, имеют хотя бы среднее школьное образование?

```
1 data.loc[data['salary'] == '>50K', 'education'].unique() # No
```



7. Отображение статистики возраста для каждой расы (функция расы) и каждого пола. Используйте `groupby()` и `describe()`. Найти максимальный возраст мужчин Амер-Индо-эскимосской расы.

```
1 for (race, sex), sub_df in data.groupby(['race', 'sex']):  
2     print("Race: {0}, sex: {1}".format(race, sex))  
3     print(sub_df['age'].describe())
```



This file was updated remotely or in another tab. To force a save, overwriting the last update, select Save from the File menu



This file was updated remotely or in another tab. To force a save, overwriting the last update, select Save from the File menu



8. Среди кого больше доля тех, кто много зарабатывает (>50 тыс.): среди женатых или одиноких мужчин (характеристика семейного положения)? Считайте женатыми тех, кто имеет семейное положение, начиная с женатого, остальные считаются холостяками.

```
1 data.loc[(data['sex'] == 'Male') &
2          (data['marital-status'].isin(['Never-married',
3                                         'Separated',
4                                         'Divorced',
5                                         'Widowed']))], 'salary'].value_counts()
```




```
1 data.loc[(data['sex'] == 'Male') &  
2          (data['marital-status'].str.startswith('Married')), 'salary'].value_count
```



```
1 data['marital-status'].value_counts()
```




9. Каково максимальное количество часов, которое человек работает в неделю (функция часов в неделю)? Сколько человек работает такое количество часов и каков процент тех, кто много зарабатывает среди них?

```
1 max_load = data['hours-per-week'].max()  
2 print("Max time - {0} hours./week.".format(max_load))  
3  
4 num_workaholics = data[data['hours-per-week'] == max_load].shape[0]  
5 print("Total number of such hard workers {0}".format(num_workaholics))  
6  
7 rich_share = float(data[(data['hours-per-week'] == max_load)  
8                       & (data['salary'] == '>50K')].shape[0]) / num_workaholics  
9 print("Percentage of rich among them {0}%".format(int(100 * rich_share)))
```



10. Подсчитайте среднее время работы (часы-в-неделю) тех, кто зарабатывает мало и много (зарплата) для каждой страны (полной страны)

This file was updated remotely or in another tab. To force a save, overwriting the last update, select Save from the File menu 

```
1 for (country, salary), sub_df in data.groupby(['native-country', 'salary']):  
2     print(country, salary, round(sub_df['hours-per-week'].mean(), 2))
```



This file was updated remotely or in another tab. To force a save, overwriting the last update, select
Save from the File menu



Элегантный метод:

```
1 pd.crosstab(data['native-country'], data['salary'],
2             values=data['hours-per-week'], aggfunc=np.mean).T
```



Часть 2

Загружаем данные

```
1 # импортируем Pandas и Numpy
2 import pandas as pd
3 import numpy as np
```

```
1 dictionary = pd.read_csv('dictionary.csv')
2 dictionary.head()
```

This file was updated remotely or in another tab. To force a save, overwriting the last update, select Save from the File menu



```
1 summer = pd.read_csv('summer.csv')
2 summer.head()
```



```
1 winter = pd.read_csv('winter.csv')  
2 winter.head()
```



Работа с данными

```
1 dictionary.head()
```



This file was updated remotely or in another tab. To force a save, overwriting the last update, select
Save from the File menu



```
1 | winter.head()
```



Соединение таблиц

```
1 | def connection_pandas(dictionary, summer):
2 |     result = pd.merge(dictionary, summer, left_on = 'Code', right_on = 'Cc
3 |     return result
4 |
5 | connection_pandas(dictionary, summer).head()
```



```
1 | import pandasql as ps
2 | pysql = lambda a: ps.sqldf(a, globals())
3 | def connection_pandasql(dictionary, summer):
4 |     query = "select * from dictionary, summer where dictionary.Code = summer.Cc
5 |     join_result = pysql(query)
6 |     return join_result
7 | abc = connection_pandasql(dictionary, summer)
8 | connection_pandasql(dictionary, summer).head()
```

This file was updated remotely or in another tab. To force a save, overwriting the last update, select Save from the File menu



Сравнение времени выполнения запросов

```
1 | import time
```

```

2 class Profiler(object):
3     def __enter__(self):
4         self._startTime = time.time()
5
6     def __exit__(self, type, value, traceback):
7         print("Elapsed time: {:.3f} sec".format(time.time() - self._startTime))
8
9 with Profiler() as p:
10     connection_pandas(dictionary, summer)

```



```

1 with Profiler() as p:
2     connection_pandas(dictionary, winter)

```



```

1 with Profiler() as p:
2     connection_pandasql(dictionary, summer)

```



```

1 with Profiler() as p:
2     connection_pandasql(dictionary, winter)

```



Вывод: соединение с помощью pandas работает в 30 быстрее, чем pandasql

Агрегирование: произвольный запрос на группировку набора данных с использованием функций агрегирования

```

1 def aggregation_pandas(dictionary, summer):
2     result = pd.merge(dictionary, summer, left_on = 'Code', right_on = 'Country_x')
3     final_0 = result[result['Year'] == 2012]
4     final = final_0[final_0['Medal'] == 'Gold'].groupby("Country_x").agg({
5         "Medal": "count",
6         "Discipline": "nunique",
7         "Gender": "nunique",
8     })
9     return final
10
11 aggregation_pandas(dictionary, summer).head(10)

```



This file was updated remotely or in another tab. To force a save, overwriting the last update, select Save from the File menu



```

1 def aggregation_pandasql(summer):
2     query = '''
3     SELECT Country, count(Medal), count(DISTINCT Discipline), count(DISTINCT (
4     WHERE Medal == 'Gold' and Year == 2012 and Country != 'None'
5     GROUP BY Country
6     '''
7     return ps.sqldf(query, locals())
8
9 aggregation_pandasql(summer).head(10)

```

	Country	count(Medal)	count(DISTINCT Discipline)	count(DISTINCT Gender)
0	ALG	1	1	1
1	ARG	1	1	1
2	AUS	19	5	2
3	AZE	2	1	1
4	BAH	4	1	1
5	BLR	3	2	2
6	BRA	14	3	2
7	CAN	1	1	1
8	CHN	56	13	2
9	COL	1	1	1

Сравнение времени выполнения запросов агрегирования

```

1 import seaborn
2 import matplotlib.pyplot as plt
3 with Profiler() as p:
4     aggregation_pandas(dictionary, summer)

```

Elapsed time: 0.025 sec

This file was updated remotely or in another tab. To force a save, overwriting the last update, select Save from the File menu

```

2 aggregation_pandasql(summer)

```

Вывод: pandas работает значительно быстрее, чем pandasql (в 10 раз)

This file was updated remotely or in another tab. To force a save, overwriting the last update, select
Save from the File menu

