

Семинар: знакомство с Panda

```
B [1]: 1 import pandas as pd
        2 import numpy as np
```

Задания для самостоятельного решения

Часть 1. Для датасета пассажиров Титаника

```
B [3]: 1 df = pd.read_csv('Data/titanic_train.csv', sep=';', encoding='utf-8')
        2 df.head()
```

```
Out[3]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	Survived
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...)	female	38.0	
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	
4	5	0	3	Allen, Mr. William Henry	male	35.0	

1. Какова доля семей, в которых минимальный возраст н детейми)?

```
B [4]: 1 df['Surname'] = df['Name'].apply(lambda name:
        2 np.sum(df.groupby('Surname')['Age'].apply(min)))
```

```
Out[4]: 0.17841079460269865
```

2. Какова доля выживших пассажиров класса 3? А пасса;

```
B [5]: 1 print("Доля выживших третьего класса: ")
      2 print(len(df[(df['Pclass'] == 3) & (df['Survived'] == 1)]))
      3 print("Доля выживших первого класса: ")
      4 print(len(df[(df['Pclass'] == 1) & (df['Survived'] == 1)]))
```

Доля выживших третьего класса

0.24236252545824846

Доля выживших первого класса

0.6296296296296297

3. Сколько пассажиров выжило, а сколько - нет?

```
B [6]: 1 print(np.sum(df['Survived'] == 1))
      2 print(np.sum(df['Survived'] == 0))
```

342

549

4. Создайте столбец 'IsChild', который равен 1, если возраст пассажира меньше 16 лет. Для пропущенных значений поведение функции может быть любым.

```

B [7]: 1 def fun(i):
        2     if i < 20:
        3         return 1
        4     else:
        5         return 0
        6
        7 df["IsChild"] = df["Age"].apply(fun)
        8 df

```

```

Out[7]:

```

	PassengerId	Survived	Pclass	Name	Sex	Age
0	1	0	3	Braund, Mr. Owen Harris	male	22.0
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0
2	3	1	3	Heikkinen, Miss. Laina	female	26.0
3	4	1	1	Futelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0
4	5	0	3	Allen, Mr. William Henry	male	35.0
...
886	887	0	2	Montvila, Rev. Juozas	male	27.0
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN
889	890	1	1	Behr, Mr. Karl Howell	male	26.0
890	891	0	3	Dooley, Mr. Patrick	male	32.0

891 rows × 14 columns



5. Какова доля выживших женщин из первого класса? А из второго класса?

```
In [8]: 1 print('Доля выживших женщин из первого класса')
2 print(len(df[(df['Pclass'] == 1) & (df['Sex']
3         / len(df[(df['Pclass'] == 1) & (df['Sex']
4 print('Доля выживших мужчин из третьего класса')
5 print(len(df[(df['Pclass'] == 3) & (df['Sex']
6         / len(df[(df['Pclass'] == 3) & (df['Sex']
```

Доля выживших женщин из первого класса

0.9680851063829787

Доля выживших мужчин из третьего класса

0.13544668587896252

Задания для самостоятельного решения

Часть 2.

Описание данных

В папке Data находится информация о студентах. Всего 10
делятся на две категории: * Students_info_i - информация с
Students_marks_i - оценки студентов из группы i за экзамен

Здесь можно посмотреть подробно о некоторых методах р:
данными pandas: [https://www.kaggle.com/residentmario/rena-combining#Combining_\(https://www.kaggle.com/residentmario/combining#Combining\)](https://www.kaggle.com/residentmario/rena-combining#Combining_(https://www.kaggle.com/residentmario/combining#Combining))

**Задача 2.1. Соберите всю информацию о студентах в о
получившейся таблице должна быть информация и оц
всех групп. Напечатайте несколько строк таблицы для
результата.¶**

```
B [9]: 1 df = pd.read_csv('Data/Students_info_0.csv').n
2 for i in range(1, 9):
3     df_1 = pd.read_csv(f'Data/Students_info_{i
4     df = pd.concat([df, df_1], ignore_index=Tr
5
6 df.head()
```

Out[9]:

	index	gender	race/ethnicity	parental level of education	lunch	prep
0	0	female	group B	bachelor's degree	standard	
1	1	female	group C	some college	standard	col
2	2	female	group B	master's degree	standard	
3	3	male	group A	associate's degree	free/reduced	
4	4	male	group C	some college	standard	

Задание 2.2. Удалите столбец index у полученной таблице 10 строк таблицы.

```
B [10]: 1 df = df.drop(['index'], axis=1)
2 df.head(10)
```

Out[10]:

	gender	race/ethnicity	parental level of education	lunch	test preparation course
0	female	group B	bachelor's degree	standard	none
1	female	group C	some college	standard	completed
2	female	group B	master's degree	standard	none
3	male	group A	associate's degree	free/reduced	none
4	male	group C	some college	standard	none
5	female	group B	associate's degree	standard	none
6	female	group B	some college	standard	completed
7	male	group B	some college	free/reduced	none
8	male	group D	high school	free/reduced	completed
9	female	group B	high school	free/reduced	none

Задание 2.3. Выведите на экран размеры полученной т

```
B [11]: 1 df.shape
```

```
Out[11]: (900, 9)
```

**Задание 2.4. Выведите на экран статистические характ
столбцов таблицы (минимум, максимум, среднее значе
отклонение)**

```
B [12]: 1 df.describe()
```

```
Out[12]:
```

	math score	reading score	writing score
count	900.000000	900.000000	900.000000
mean	66.045556	69.085556	67.910000
std	15.099269	14.568942	15.213194
min	0.000000	17.000000	10.000000
25%	57.000000	59.000000	57.000000
50%	66.000000	70.000000	69.000000
75%	77.000000	79.000000	79.000000
max	100.000000	100.000000	100.000000

Задание 2.5. Проверьте, есть ли в таблице пропущенны

```
B [13]: 1 df.isnull().sum()
```

```
Out[13]: gender                                0
race/ethnicity                                0
parental level of education                    0
lunch                                           0
test preparation course                         0
group                                           0
math score                                     0
reading score                                 0
writing score                                 0
dtype: int64
```

**Задание 2.6. Выведите на экран средние баллы студен
(math, reading, writing)**

```
B [14]: 1 df.describe().iloc[1]
```

```
Out[14]: math score      66.045556
reading score    69.085556
writing score    67.910000
Name: mean, dtype: float64
```

**Задание 2.7. Как зависят оценки от того, проходил ли с
подготовки к сдаче экзамена (test preparation course)?**

каждого предмета в отдельности средний балл студента для подготовки к экзамену и не проходивших курс.

```
B [15]: 1 print("Средние баллы учеников, проходивших курс  
2 df[df['test preparation course'] == 'completed']
```

Средние баллы учеников, проходивших курс подготовки

```
Out[15]: math score      69.814465  
reading score    74.091195  
writing score    74.613208  
dtype: float64
```

```
B [16]: 1 print('Средние баллы учеников, не проходивших  
2 df[df['test preparation course'] == 'none'].agg
```

Средние баллы учеников, не проходивших курс подготовки

```
Out[16]: math score      63.986254  
reading score    66.350515  
writing score    64.247423  
dtype: float64
```

Задание 2.8. Выведите на экран все различные значения

```
B [17]: 1 [print(i) for i in df.lunch.unique()]
```

standard
free/reduced

```
Out[17]: [None, None]
```

Задание 2.9. Переименуйте колонку "parental level of education" в "test preparation" с помощью rename

<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.rename.html>
(<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.rename.html>)

```
B [18]: 1 df = df.rename(columns={'parental level of edu
2 df.head()
```

```
Out[18]:
```

	gender	race/ethnicity	education	lunch	test preparation
0	female	group B	bachelor's degree	standard	none
1	female	group C	some college	standard	completed
2	female	group B	master's degree	standard	none
3	male	group A	associate's degree	free/reduced	none
4	male	group C	some college	standard	none

Зафиксируем минимальный балл для сдачи экзамена

passmark = 50

Задание 2.10. Ответьте на вопросы:

- * Какая доля студентов сдала экзамен по математике (passmark)?
- * Какая доля студентов, проходивших курс подготовки к экзамену по математике?
- * Какая доля девушек, не проходивших курс подготовки к экзамену по математике?

```
B [19]: 1 len(df[df['math score'] >= 50]) / len(df)
```

```
Out[19]: 0.8666666666666667
```

```
B [20]: 1 len(df[(df['math score'] >= 50) & (df['test preparation'] == 'completed')]) / len(df)
2
```

```
Out[20]: 0.9245283018867925
```

```
B [21]: 1 len(df[(df['math score'] >= 50) & (df['test preparation'] == 'completed') & (df['gender'] == 'female')]) / len(df)
2
```

```
Out[21]: 0.9156626506024096
```

Задание 2.11. С помощью groupby выполните задания по времени выполнения каждого из заданий.

- * Для каждой этнической группы выведите средний балл за экзамен по математике.
- * Для каждого уровня образования выведите минимальный балл за экзамен по математике.


```
B [22]: 1 df1 = df.groupby(['race/ethnicity']).agg({'reading score': 'mean'})
        2 df1
```

Out[22]:

reading score	
race/ethnicity	
group A	65.194805
group B	67.359551
group C	69.037801
group D	69.643478
group E	73.056452

```
B [23]: 1 %%timeit
        2 df.groupby(['race/ethnicity']).agg({'reading score': 'mean'})
```

536 μ s \pm 10.4 μ s per loop (mean \pm std. dev. of 7 runs)

```
B [24]: 1 df2 = df.groupby(['education']).agg({'writing score': 'mean'})
        2 df2
```

Out[24]:

writing score	
education	
associate's degree	35
bachelor's degree	38
high school	15
master's degree	46
some college	19
some high school	10

```
B [25]: 1 %%timeit
        2 df.groupby(['education']).agg({'writing score': 'mean'})
```

594 μ s \pm 27.1 μ s per loop (mean \pm std. dev. of 7 runs)

Задание 2.12. Выполните задание 11 с помощью цикла выполнения.

```
B [26]: 1 df1_1 = pd.DataFrame({'race/ethnicity': df1.index, 'reading score': df1['reading score']})
        2 s = df[['race/ethnicity', 'reading score']]
```

```
B [27]: 1 %%timeit
        2 for i in df1_1.index:
        3     df1_1.loc[i] = (s[s['race/ethnicity'] == df1_1.index[i]]['reading score'].values[0])
        4 df1_1
```

1.33 ms \pm 5.89 μ s per loop (mean \pm std. dev. of 7 runs)

```

B [28]: 1 df2_2 = pd.DataFrame({'education': df2.index,
2 s = df[['education', 'writing score']]

B [29]: 1 %%timeit
2 for i in df2_2.index:
3     df2_2.loc[i] = min((s[s['education'] == i]
4 df2_2

```

1.38 ms ± 8.66 µs per loop (mean ± std. dev. of 7

Выполнение задания с помощью циклов занимает больше этих же заданий с помощью groupby

Задание 2.13. Выведите на экран средние баллы студента по предмету в зависимости от пола и уровня образования. Для каждого уровня образования получите количество групп, равных 2 * (число уровней образования). Для каждой такой группы выведите средний балл по каждому предмету.

Это можно сделать с помощью сводных таблиц (pivot_table)

<https://www.kaggle.com/kamilpolak/tutorial-how-to-use-pivot-table>
(<https://www.kaggle.com/kamilpolak/tutorial-how-to-use-pivot-table>)

```

B [30]: 1 df.groupby(['gender', 'education']).agg({'math score': 'mean', 'reading score': 'mean', 'writing score': 'mean'})

```

Out[30]:

		math score	reading score	writing score
gender	education			
female	associate's degree	65.457944	74.093458	74.1
	bachelor's degree	68.232143	77.053571	78.2
	high school	61.012048	69.626506	68.2
	master's degree	64.468750	74.906250	75.9
	some college	65.446602	73.417476	73.8
	some high school	59.168675	69.192771	68.3
male	associate's degree	71.020202	67.797980	65.6
	bachelor's degree	70.431373	67.921569	67.2
	high school	63.806818	61.045455	57.7
	master's degree	74.476190	72.952381	71.8
	some college	68.479592	64.306122	62.4
	some high school	67.392405	64.189873	60.5

Задание 2.14. Сколько студентов успешно сдали экзамен по математике?

Создайте новый столбец в таблице df под названием Math_F, если студент не сдал экзамен по математике (балл меньше 60), иначе.

Посчитайте количество студентов, сдавших и не сдавших экзамен по математике.

Сделайте аналогичные шаги для экзаменов по чтению и пи

```
B [31]: 1 df.loc[df['math score'] >= 50, 'Math_PassStatu
2 df.loc[df['math score'] < 50, 'Math_PassStatus
3
4 df.loc[df['reading score'] >= 50, 'Reading_Pas
5 df.loc[df['reading score'] < 50, 'Reading_Pass
6
7 df.loc[df['writing score'] >= 50, 'Writing_Pas
8 df.loc[df['writing score'] < 50, 'Writing_Pass
9
10 df
```

```
Out[31]:
```

	gender	race/ethnicity	education	lunch	te preparatic
0	female	group B	bachelor's degree	standard	no
1	female	group C	some college	standard	complete
2	female	group B	master's degree	standard	no
3	male	group A	associate's degree	free/reduced	no
4	male	group C	some college	standard	no
...
895	female	group E	some high school	free/reduced	no
896	male	group B	high school	free/reduced	no
897	female	group B	some high school	free/reduced	complete
898	male	group D	associate's degree	standard	complete
899	female	group D	some high school	standard	complete

900 rows × 12 columns

Задание 2.15. Сколько студентов успешно сдали все эк

Создайте столбец OverAll_PassStatus и запишите в него др
если студент не сдал хотя бы один из трех экзаменов, а ин

Посчитайте количество студентов, которые сдали все экза

B [36]:

1

2

3

df.loc[(df['math score'] >= 50) & (df['reading score'] >= 50)]

df.loc[(df['math score'] < 50) | (df['reading score'] < 50)]

df

Out[36]:

	gender	race/ethnicity	education	lunch	test preparation course
0	female	group B	bachelor's degree	standard	not completed
1	female	group C	some college	standard	completed
2	female	group B	master's degree	standard	not completed
3	male	group A	associate's degree	free/reduced	not completed
4	male	group C	some college	standard	not completed
...
895	female	group E	some high school	free/reduced	not completed
896	male	group B	high school	free/reduced	not completed
897	female	group B	some high school	free/reduced	completed
898	male	group D	associate's degree	standard	completed
899	female	group D	some high school	standard	completed

900 rows × 6 columns