

Библиотека Pandas

[Материалы](#) > [Анализ и обработка данных](#)

Мы начинаем курс анализа и обработки данных. На двух предыдущих курсах мы, во-первых, получили общее представление о том, как устроено машинное обучение, и, во-вторых, приобрели достаточно продвинутые навыки в [программировании на Питоне](#).

Курс анализа данных — это первый курс, в рамках которого мы переходим непосредственно к работе над построением моделей.

Для того чтобы лучше понимать содержание этого курса, давайте вначале рассмотрим этапы построения модели.

Задача машинного обучения

Ниже представлена общая последовательность работы над задачей машинного обучения.



Описанный процесс принято называть **пайплайном**, то есть порядком действий (от англ. pipeline, трубопровод, конвейер), которые необходимо выполнить для построения модели. Подробнее рассмотрим некоторые из этих этапов.

Этап 1. Постановка задачи и определение метрики

Первый этап может показаться тривиальным, однако во многих случаях, особенно в задачах классификации, **выбор правильной метрики** является ключевым для построения качественной модели. Про важность классификационной метрики мы начали говорить в рамках вводного курса и продолжим этот разговор в дальнейшем на гораздо более детальном уровне.

Этап 2. Получение данных

Важный этап. Хотя на этом курсе мы будем использовать уже готовые (зачастую учебные) датасеты, стоит помнить, что получение данных (data gathering) не происходит само собой и во многом от того как и какие данные собраны будет зависеть конечный результат (на который не смогут повлиять ни качественный EDA, ни сложный алгоритм машинного обучения).

Этап 3. Исследовательский анализ данных

В рамках EDA нам нужно решить три основные задачи: описать данные, найти отличия и выявить взаимосвязи. **Описание данных** позволяет понять, о каких данных идет речь, а также выявить недостатки в данных (с которыми мы справляемся на этапе обработки).

Отличия и взаимосвязи в данных — основа для построения модели, это то, за что модель «цепляется», чтобы выдать верный числовой результат, правильно классифицировать или сформировать кластер.

Для решения этих задач наилучшим образом подходят средства визуализации и описательная статистика. Этим мы займемся во втором разделе.

Отдельно хочется сказать про **baseline models**, простые модели, которые мы строим в самом начале работы. Они позволяют понять, какой результат мы можем получить, не вкладывая дополнительных усилий в работу с данными, а затем отталкиваться от этого результата для обработки данных и построения более сложных моделей.

Базовые модели мы начнем строить на курсе по обучению модели.

Этап 4. Обработка данных

Как уже было сказано, на этапе EDA зачастую становится очевидно, что в данных есть недостатки, которые сильно повлияют на качество модели или в целом не позволят ее обучить.

Очистка данных: ошибки и пропуски

Во-первых, в данных могут встречаться **ошибки**: дубликаты, неверные значения или неподходящий формат данных. Кроме того, данные могут содержать **пропуски**, и с ними также нужно что-то делать. Этим вопросам посвящен третий раздел курса.

Преобразование данных

Во-вторых, зачастую *количественные данные* нужно привести к одному масштабу и/или нормальному распределению. Кроме того, числовые признаки могут содержать сильно отличающиеся от остальных данных значения или выбросы, которые также повлияют на конечный результат. *Категориальные данные* необходимо закодировать с помощью чисел. Если категориальные данные выражены строками, это может воспрепятствовать обучению алгоритма.

Преобразование количественных и категориальных данных рассматривается в четвертом разделе.

Конструирование и отбор признаков, понижение размерности

Еще одним важным этапом является *конструирование признаков*, а также *отбор признаков* и *понижение размерности*. В рамках этого курса мы затронем лишь базовые способы конструирования признаков. Более сложные вопросы отбора признаков и понижения размерности мы отложим на потом.

Этап 5. Моделирование и оценка результата

Когда данные готовы, их можно загружать в модель, обучать эту модель и оценивать результат.

Здесь важно сказать, что это *итеративный* (iterative) или циклический процесс. Многие из описанных выше шагов могут повторяться. В частности, построение модели может привести к необходимости дополнительной обработки данных. Кроме того, разные алгоритмы требуют разной подготовки данных (например, линейные модели требуют масштабирования данных, а для деревьев решений этого не нужно).

При этом прежде чем приступить к анализу и обработке данных, важно освоить библиотеку Pandas. Именно этим мы и займемся в начале курса.

Про библиотеку Pandas

Библиотека Pandas — это ключевой инструмент для анализа данных в Питоне. Она позволяет работать с данными, представленными в табличной форме, а также временными рядами. Как вы уже видели, Pandas легко интегрируется с matplotlib, seaborn, sklearn и другими библиотеками.

Кроме того, структурно изучение Pandas можно разделить на две части: преобразование данных и статистический анализ. В этом разделе (первое и второе занятие) мы начнем знакомиться с первой частью, а в следующем (занятия три и четыре) перейдем ко второй.

Откроем блокнот к этому занятию📓

Объекты DataFrame и Series

В Pandas есть два основных объекта: DataFrame и Series. Слегка упрощая, можно сказать, что DataFrame — это таблица (и соответственно двумерный массив), а Series — столбец этой таблицы (а значит одномерный массив).

The diagram shows a DataFrame with 5 columns and 5 rows. A red box highlights the entire table. A red arrow points to the header row with the label 'columns'. A red arrow points to the first column with the label 'index'. A red arrow points to the first row with the label 'values'. A green arrow points to the fourth column with the label 'Series'.

	Column 1	Column 2	Column 3	Column 4	Column 5
0	Value 1.1	Value 2.1	Value 3.1	Value 4.1	Value 5.1
1	Value 1.2	Value 2.2	Value 3.2	Value 4.2	Value 5.2
2	Value 1.3	Value 2.3	Value 3.3	Value 4.3	Value 5.3
3	Value 1.4	Value 2.4	Value 3.4	Value 4.4	Value 5.4
4	Value 1.5	Value 2.5	Value 3.5	Value 4.5	Value 5.5

Индекс (index), столбцы (columns) и значения (values) датафрейма мы разберем чуть позднее, а сейчас давайте посмотрим на способы создания датафрейма.

Создание датафрейма

Способ 1. Создание датафрейма из файла

Датафрейм можно создать, импортировав файлы различных форматов. Мы уже делали так с файлами в формате .csv с помощью функции `pd.read_csv()`. Аналогичным образом можно импортировать файлы и в других форматах, например MS Excel или html. Рассмотрим несколько примеров.

Файл .csv в zip-архиве. Если в zip-архиве содержится только один файл, его можно напрямую подгрузить в Pandas. Вначале скачаем необходимые данные.

train.zip

[Скачать](#)

После скачивания не забудьте подгрузить данные в сессионное хранилище Google Colab.

```
1 # импортируем файл из папки content и выведем первые три строки
2 csv_zip = pd.read_csv('/content/train.zip')
3 csv_zip.head(3)
```

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C

2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
---	---	---	---	------------------------	--------	------	---	---	---------------------	--------	-----	---

Как вы наверное узнали, речь идет о датасете «Титаник».

MS Excel. Похожим образом мы можем импортировать файлы в формате Excel. Скачаем файл.

iris.xlsx

[Скачать](#)

Теперь используем функцию **pd.read_excel()** для импорта файла.

```
1 # импортируем данные в формате Excel, указав номер листа, который хотим использовать
2 excel_data = pd.read_excel('/content/iris.xlsx', sheet_name = 0)
3 excel_data.head(3)
```

	Petal_width	Petal_length	Sepal_width	Sepal_length	Species_name
0	0.2	1.4	3.5	5.1	Setosa
1	0.2	1.4	3.0	4.9	Setosa
2	0.2	1.3	3.2	4.7	Setosa

С [датасетом про цветы ириса](#) мы также уже знакомы.

Дополнительно запомним, что файлы можно хранить локально на своем компьютере (если вы используете [Jupyter Notebook](#)), импортировать в сессионное хранилище Google Colab или передать в функцию импорта ссылку на файл в Интернете.

Файл в формате html. Данные можно загружать не только из файла, но и брать из Интернета. Для этого подойдет **функция read_html()**. По большому счету речь идет о примере несложного **веб-скрапинга** (web scraping), то есть получения информации с веб-страниц в автоматизированном режиме.

Возьмем страницу Википедии про [население Земли](#)^[6]. В частности предположим, что нас будет интересовать вот эта таблица.

World population milestones in billions ^[6] (Worldometers estimates)										
Population	1	2	3	4	5	6	7	8	9	10
Year	1804	1927	1960	1974	1987	1999	2011	2022	2037	2057
Years elapsed	200,000+	123	33	14	13	12	12	11	15	20

```
1 # передадим соответствующую ссылку в функцию pd.read_html()
2 # в параметре match укажем ключевые слова, которые помогут найти нужную таблицу
3 html_data = pd.read_html('https://en.wikipedia.org/wiki/World_population',
4                           match = 'World population')
5
6 # мы получили пять результатов
7 len(html_data)
```

```
1 5
```

Посмотрим на первый полученный результат.

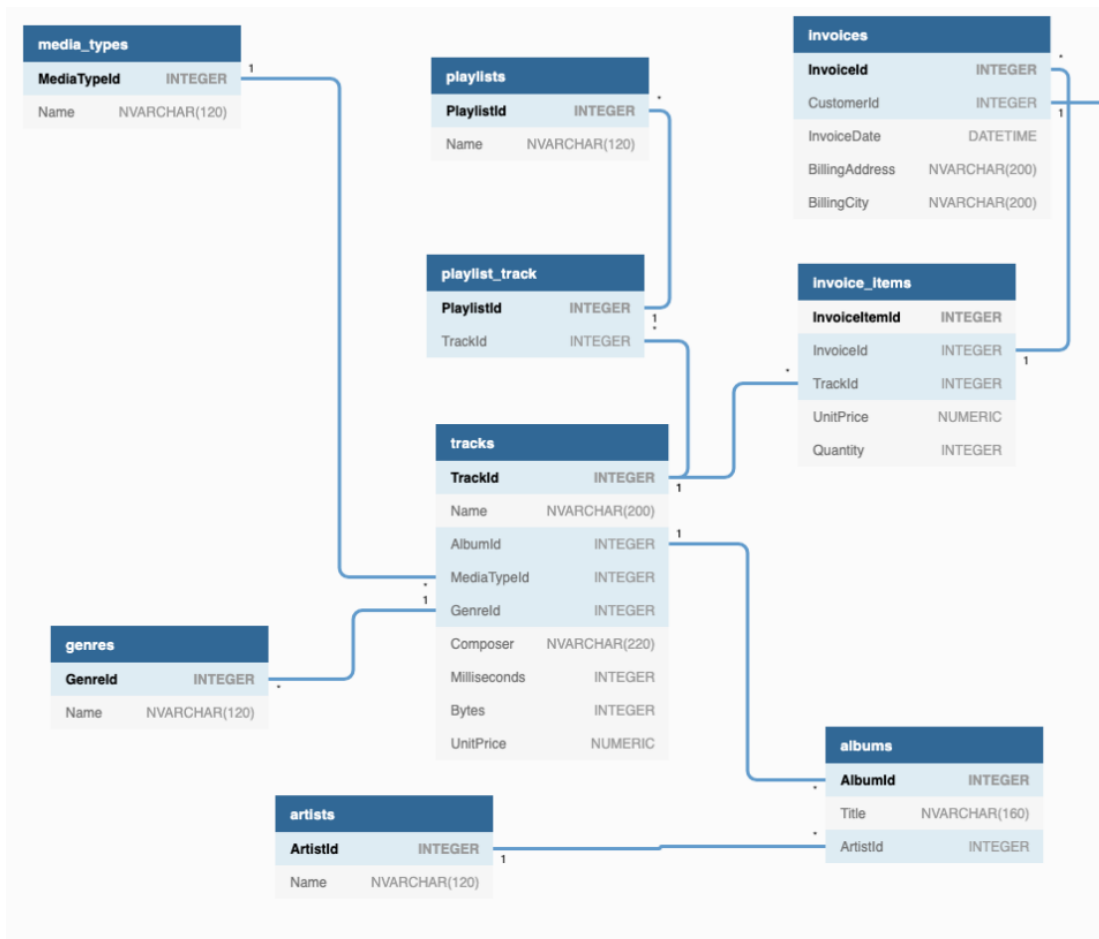
```
1 | html_data[0]
```

	Population	1	2	3	4	5	6	7	8	9	10
0	Year	1804	1927	1960	1974	1987	1999	2011	2022	2037	2057
1	Years elapsed	200,000+	123	33	14	13	12	12	11	15	20

Способ 2. Подключение к базе данных SQL

Библиотека Pandas позволяет получать информацию напрямую из базы данных SQL. В качестве примера возьмем популярную учебную базу данных Chinook, в которой содержится информация о магазине, торгующем цифровыми записями (digital media store).

Ниже представлена схема этой **реляционной базы данных** (relational database schema).



Тема реляционных баз данных и SQL выходит за рамки сегодняшнего занятия, однако для полноты изложения скажу, что такая база данных называется **реляционной** (от англ. relation — «отношение, зависимость»), потому что ее таблицы взаимосвязаны.

Например, в таблице *tracks* (записи) в поле *GenreId* содержится идентификатор жанра. По этому идентификатору мы можем определить, к какому жанру принадлежит та или иная композиция, заглянув в таблицу *genres*, где каждому *GenreId* соответствует название определенного жанра.

SQL же расшифровывается как structured query language (**язык структурированных запросов**) и позволяет достаточно гибко создавать, извлекать и изменять данные из реляционных баз данных.

Попробуем подключить базу данных к нашему блокноту и написать очень простой запрос на SQL внутри кода на Питоне. Вначале скачаем базу данных chinook (архив нужно распаковать и подгрузить файл в формате .db).

[chinook.zip](#)[Скачать](#)

```
1 # импортируем модуль sqlite3 для работы с базой данных SQL
2 import sqlite3 as sql
3
4 # создадим соединение с базой данных chinook
5 conn = sql.connect('/content/chinook.db')
6
7 # выберем все строки из таблицы tracks
8 sql_data = pd.read_sql('SELECT * FROM tracks', conn) # vs. read_sql_query
9
10 # посмотрим на результат
11 sql_data.head(3)
```

	TrackId	Name	AlbumId	MediaTypeId	GenreId	Composer	Milliseconds	Bytes	UnitPrice
0	1	For Those About To Rock (We Salute You)	1	1	1	Angus Young, Malcolm Young, Brian Johnson	343719	11170334	0.99
1	2	Balls to the Wall	2	2	1	None	342562	5510424	0.99
2	3	Fast As a Shark	3	2	1	F. Baltes, S. Kaufman, U. Dirkschneider & W. Ho...	230619	3990994	0.99

Позднее в рамках курса мы поговорим еще про два формата: json и pickle.

Способ 3. Создание датафрейма из словаря

Кроме этого, и мы уже много раз так делали, библиотека Pandas позволяет создать датафрейм из словаря с помощью функции `pd.DataFrame()`. В этом случае ключи словаря будут названиями столбцов, а значения — элементами строк. Последние можно передать с помощью списков или массивов Numpy. Рассмотрим на примере.

Вначале подготовим данные.

```
1 # создадим несколько списков и массивов Numpy с информацией о семи странах мира
2 country = np.array(['China', 'Vietnam', 'United Kingdom', 'Russia', 'Argentina', 'Bolivia'])
3 capital = ['Beijing', 'Hanoi', 'London', 'Moscow', 'Buenos Aires', 'Sucre', 'Pretoria']
4 population = [1400, 97, 67, 144, 45, 12, 59] # млн. человек
5 area = [9.6, 0.3, 0.2, 17.1, 2.8, 1.1, 1.2] # млн. кв. км.
6 sea = [1] * 5 + [0, 1] # выход к морю (в этом списке его нет только у Боливии)
```

Теперь создадим пустой словарь и поместим в него наши данные.

```
1 # создадим пустой словарь
2 countries_dict = {}
3
4 # превратим эти списки в значения словаря,
```

```
5 # одновременно снабдив необходимыми ключами
6 countries_dict['country'] = country
7 countries_dict['capital'] = capital
8 countries_dict['population'] = population
9 countries_dict['area'] = area
10 countries_dict['sea'] = sea
```

Наконец, создадим датафрейм с помощью функции `pd.DataFrame()`.

```
1 countries = pd.DataFrame(countries_dict)
2 countries
```

	country	capital	population	area	sea
0	China	Beijing	1400	9.6	1
1	Vietnam	Hanoi	97	0.3	1
2	United Kingdom	London	67	0.2	1
3	Russia	Moscow	144	17.1	1
4	Argentina	Buenos Aires	45	2.8	1
5	Bolivia	Sucre	12	1.1	0
6	South Africa	Pretoria	59	1.2	1

Способ 4. Создание датафрейма из 2D массива Numpy

Помимо этого, датафрейм можно создать из двумерного массива Numpy.

```
1 # внешнее измерение будет столбцами, внутреннее - строками
2 arr = np.array([[1, 1, 1],
3                 [2, 2, 2],
4                 [3, 3, 3]])
5
6 pd.DataFrame(arr)
```

	0	1	2
0	1	1	1
1	2	2	2
2	3	3	3

Перейдем к структуре и свойствам датафрейма.

Структура и свойства датафрейма

Датафрейм библиотеки Pandas состоит из трех основных компонентов: строк (index), столбцов (columns) и значений (values).

Получить доступ к названиям столбцов можно с помощью атрибута `columns`.

```
1 # выведем столбцы датафрейма countries
2 countries.columns
```

```
1 Index(['country', 'capital', 'population', 'area', 'sea'], dtype='object')
```


Атрибут index описывает, каким образом идентифицируются строки.

```
1 # в нашем случае индекс это числовая последовательность
2 countries.index
```

```
1 RangeIndex(start=0, stop=7, step=1)
```

С помощью **атрибута values** мы получаем доступ к значениям датафрейма.

```
1 # значения выводятся в формате массива Numpy
2 countries.values
```

```
1 array([[ 'China', 'Beijing', 1400, 9.6, 1],
2        [ 'Vietnam', 'Hanoi', 97, 0.3, 1],
3        [ 'United Kingdom', 'London', 67, 0.2, 1],
4        [ 'Russia', 'Moscow', 144, 17.1, 1],
5        [ 'Argentina', 'Buenos Aires', 45, 2.8, 1],
6        [ 'Bolivia', 'Sucre', 12, 1.1, 0],
7        [ 'South Africa', 'Pretoria', 59, 1.2, 1]], dtype=object)
```

Описание индекса и названия столбцов мы также можем получить с помощью **атрибута axes** (множественное от axis, ось) с индексами [0] и [1] соответственно.

```
1 # выведем описание индекса датафрейма через атрибут axes[0]
2 countries.axes[0]
```

```
1 RangeIndex(start=0, stop=7, step=1)
```

```
1 # axes[1] выводит названия столбцов
2 countries.axes[1]
```

```
1 Index([ 'country', 'capital', 'population', 'area', 'sea'], dtype='object')
```

В целом, полезно запомнить, что axis = 0 позволяет применить какой-либо метод к строкам датафрейма, а axis = 1, соответственно к столбцам.

Мы также можем посмотреть количество измерений, размерность и общее количество элементов датафрейма с помощью **атрибутов ndim, shape и size** соответственно.

```
1 countries.ndim, countries.shape, countries.size
```

```
1 (2, (7, 5), 35)
```

Атрибут dtypes выдает тип данных каждого столбца.

```
1 countries.dtypes
```

```
1 country      object
2 capital      object
3 population    int64
4 area         float64
5 sea          int64
6 dtype: object
```

Кроме того, с помощью **метода .memory_usage()** мы можем посмотреть объем занимаемой памяти по столбцам в байтах.

```
1 countries.memory_usage()
```

```
1 Index      128
2 country     56
3 capital     56
4 population  56
```

```
5 | area          56
6 | sea           56
7 | dtype: int64
```

Индекс датафрейма

По умолчанию при создании датафрейма библиотека Pandas снабжает его числовым индексом, который начинается с нуля. При этом мы можем присвоить датафрейму собственный индекс.

Присвоение индекса

Более того, этот индекс не обязательно должен быть числовым.

```
1 | # создадим список с кодами стран
2 | custom_index = ['CN', 'VN', 'GB', 'RU', 'AR', 'BO', 'ZA']
3 |
4 | # создадим датафрейм из словаря и передадим список через параметр index
5 | countries = pd.DataFrame(countries_dict,
6 |                           index = custom_index)
7 |
8 | # посмотрим на результат
9 | countries
```

	country	capital	population	area	sea
CN	China	Beijing	1400	9.6	1
VN	Vietnam	Hanoi	97	0.3	1
GB	United Kingdom	London	67	0.2	1
RU	Russia	Moscow	144	17.1	1
AR	Argentina	Buenos Aires	45	2.8	1
BO	Bolivia	Sucre	12	1.1	0
ZA	South Africa	Pretoria	59	1.2	1

Этот индекс можно сбросить с помощью метода **`.reset_index()`**.

```
1 | # при этом параметр inplace = True сделает изменения постоянными
2 | countries.reset_index(inplace = True)
3 | countries
```

	index	country	capital	population	area	sea
0	CN	China	Beijing	1400	9.6	1
1	VN	Vietnam	Hanoi	97	0.3	1
2	GB	United Kingdom	London	67	0.2	1
3	RU	Russia	Moscow	144	17.1	1
4	AR	Argentina	Buenos Aires	45	2.8	1
5	BO	Bolivia	Sucre	12	1.1	0
6	ZA	South Africa	Pretoria	59	1.2	1

Как мы видим, прошлый индекс стал отдельным столбцом. Снова сделаем этот столбец индексом с помощью метода **`.set_index()`**.

```
1 # передадим методу название столбца, который хотим сделать индексом
2 countries.set_index('index', inplace = True)
3 countries
```

	country	capital	population	area	sea
index					
CN	China	Beijing	1400	9.6	1
VN	Vietnam	Hanoi	97	0.3	1
GB	United Kingdom	London	67	0.2	1
RU	Russia	Moscow	144	17.1	1
AR	Argentina	Buenos Aires	45	2.8	1
BO	Bolivia	Sucre	12	1.1	0
ZA	South Africa	Pretoria	59	1.2	1

Еще раз сбросим индекс, но на этот раз не будем делать его отдельным столбцом. Для этого передадим методу `.reset_index()` параметр `drop = True`.

```
1 countries.reset_index(drop = True, inplace = True)
2 countries
```

	country	capital	population	area	sea
0	China	Beijing	1400	9.6	1
1	Vietnam	Hanoi	97	0.3	1
2	United Kingdom	London	67	0.2	1
3	Russia	Moscow	144	17.1	1
4	Argentina	Buenos Aires	45	2.8	1
5	Bolivia	Sucre	12	1.1	0
6	South Africa	Pretoria	59	1.2	1

Собственный индекс можно создать, просто поместив новые значения в атрибут `index`.

```
1 countries.index = custom_index
2 countries
```

	country	capital	population	area	sea
CN	China	Beijing	1400	9.6	1
VN	Vietnam	Hanoi	97	0.3	1
GB	United Kingdom	London	67	0.2	1
RU	Russia	Moscow	144	17.1	1
AR	Argentina	Buenos Aires	45	2.8	1
BO	Bolivia	Sucre	12	1.1	0
ZA	South Africa	Pretoria	59	1.2	1

Многоуровневый индекс

Многоуровневый (MultiIndex) или **иерархический** (hierarchical) индекс позволяет задать несколько уровней (levels) для индексации строк или столбцов. «Внешний» индекс

объединяет несколько элементов «внутреннего» индекса.

Вначале подготовим данные для многоуровневого индекса строк.

```
1 # создадим список из кортежей с названием континента и кодом страны
2 rows = [('Asia', 'CN'),
3         ('Asia', 'VN'),
4         ('Europe', 'GB'),
5         ('Europe', 'RU'),
6         ('S. America', 'AR'),
7         ('S. America', 'BO'),
8         ('Africa', 'ZA')]
```

Продедаем аналогичную работу для столбцов.

```
1 # в столбцах название страны и столицы мы объединим в категорию names
2 # а размер населения, площадь и выход к морю в data
3 cols = [('names', 'country'),
4         ('names', 'capital'),
5         ('data', 'population'),
6         ('data', 'area'),
7         ('data', 'sea')]
```

Теперь создадим иерархический индекс для строк и столбцов с помощью функции **pd.MultiIndex.from_tuples()**.

```
1 # создадим многоуровневый индекс для строк
2 # индексам присвоим названия через names = ['region', 'code']
3 custom_multindex = pd.MultiIndex.from_tuples(rows, names = ['region', 'code'])
4
5 # сделаем то же самое для столбцов
6 custom_multicols = pd.MultiIndex.from_tuples(cols)
```

Наконец остается передать эти индексы в атрибуты **index** и **columns** и вывести результат.

```
1 countries.index = custom_multindex
2 countries.columns = custom_multicols
3
4 countries
```

names		data				
		country	capital	population	area	sea
region	code					
Asia	CN	China	Beijing	1400	9.6	1
	VN	Vietnam	Hanoi	97	0.3	1
Europe	GB	United Kingdom	London	67	0.2	1
	RU	Russia	Moscow	144	17.1	1
S. America	AR	Argentina	Buenos Aires	45	2.8	1
	BO	Bolivia	Sucre	12	1.1	0
Africa	ZA	South Africa	Pretoria	59	1.2	1

Как мы видим, страны (code) теперь разбиты на континенты (region), а информация о них — на текстовые (names) и числовые (data) данные.

```
1 # вернемся к обычному индексу и названиям столбцов
2 custom_cols = ['country', 'capital', 'population', 'area', 'sea']
3
4 countries.index = custom_index
```

```
5 | countries.columns = custom_cols
6 |
7 | countries
```

	country	capital	population	area	sea
CN	China	Beijing	1400	9.6	1
VN	Vietnam	Hanoi	97	0.3	1
GB	United Kingdom	London	67	0.2	1
RU	Russia	Moscow	144	17.1	1
AR	Argentina	Buenos Aires	45	2.8	1
BO	Bolivia	Sucre	12	1.1	0
ZA	South Africa	Pretoria	59	1.2	1

Преобразование в другие форматы

Получившийся датафрейм можно преобразовать в словарь с помощью **метода .to_dict()**.

```
1 | print(countries.to_dict())
```

```
1 | {'country': {'CN': 'China', 'VN': 'Vietnam', 'GB': 'United Kingdom', 'RU': 'Russia', 'A
```

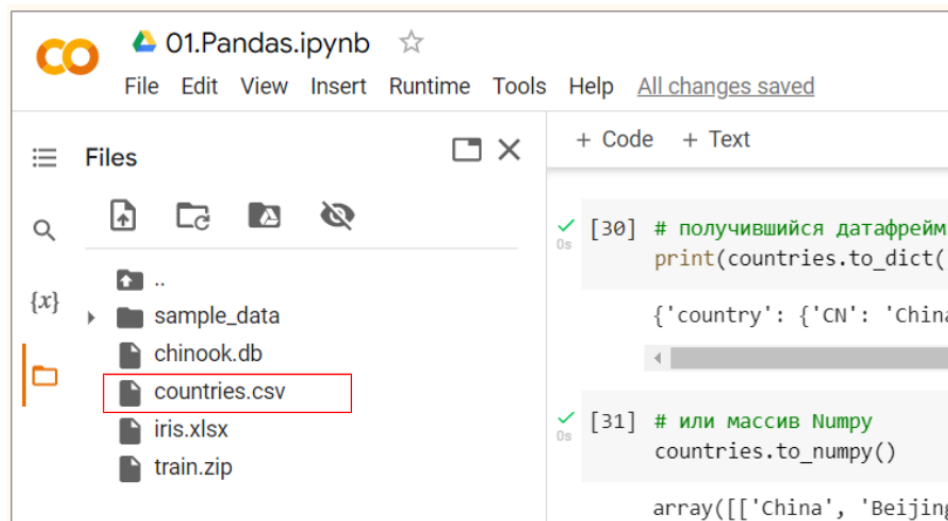
Аналогично, **метод .to_numpy()** преобразовывает данные в массив Numpy.

```
1 | countries.to_numpy()
```

```
1 | array([[ 'China', 'Beijing', 1400, 9.6, 1],
2 |        [ 'Vietnam', 'Hanoi', 97, 0.3, 1],
3 |        [ 'United Kingdom', 'London', 67, 0.2, 1],
4 |        [ 'Russia', 'Moscow', 144, 17.1, 1],
5 |        [ 'Argentina', 'Buenos Aires', 45, 2.8, 1],
6 |        [ 'Bolivia', 'Sucre', 12, 1.1, 0],
7 |        [ 'South Africa', 'Pretoria', 59, 1.2, 1]], dtype=object)
```

Кроме того, мы можем сохранить датафрейм в файл (в Google Colab он появится в сессионном хранилище). Например, это может быть файл в формате .csv и тогда необходимо использовать **функцию .to_csv()**.

```
1 | # по умолчанию, индекс также станет частью .csv файла
2 | # параметр index = False позволит этого избежать
3 | countries.to_csv('countries.csv', index = False)
```



	['Vietnam', 'Hanoi']
	['United Kingdom', 'London']
	['Russia', 'Moscow']
	['Argentina', 'Buenos Aires']
	['Bolivia', 'Sucre']
	['South Africa', 'Cape Town']

Метод `.to_list()` позволяет преобразовать объект Series (столбец датафрейма) в список.

```
1 print(countries.country.to_list())
```

```
1 ['China', 'Vietnam', 'United Kingdom', 'Russia', 'Argentina', 'Bolivia', 'South Africa']
```

Применить метод `.to_list()` к датафрейму не получится.

Создание Series

Создание Series из списка

Объект Series можно создать из списка. Возьмем следующий список с названиями стран.

```
1 country_list = ['China', 'South Africa', 'United Kingdom', 'Russia', 'Argentina', 'Vietnam']
```

Передадим его в функцию `pd.Series()`.

```
1 country_Series = pd.Series(country_list)
2 country_Series
```

```
1 0      China
2 1  South Africa
3 2  United Kingdom
4 3      Russia
5 4    Argentina
6 5      Vietnam
7 6    Australia
8 dtype: object
```

Обратите внимание, что для объекта Series был автоматически создан числовой индекс. С его помощью мы можем получить доступ к элементам этого объекта.

```
1 # например, выведем первый элемент
2 country_Series[0]
```

```
1 'China'
```

Создание Series из словаря

Помимо списка, объект Series можно создать из словаря. Создадим словарь, в котором коды стран будут ключами, а их названия — значениями.

```
1 country_dict = {'CN' : 'China',
2                 'ZA' : 'South Africa',
3                 'GB' : 'United Kingdom',
4                 'RU' : 'Russia',
5                 'AR' : 'Argentina',
6                 'VN' : 'Vietnam',
7                 'AU' : 'Australia'
8                 }
```

Передадим этот словарь в функцию **pd.Series()**. Ключи станут индексами, а значения — элементами объекта.

```
1 country_Series = pd.Series(country_dict)
2 country_Series
```

```
1 CN          China
2 ZA      South Africa
3 GB  United Kingdom
4 RU          Russia
5 AR      Argentina
6 VN      Vietnam
7 AU      Australia
8 dtype: object
```

Доступ к элементам будет осуществляться теперь уже по коду отдельной страны.

```
1 country_Series['AU']
```

```
1 'Australia'
```

Доступ к строкам, столбцам и элементам

Поговорим про способы обращения к строкам, столбцам и элементам датафрейма.

Циклы в датафрейме

Цикл **for** позволяет получить доступ к названиям столбцов датафрейма.

```
1 for column in countries:
2     print(column)
```

```
1 country
2 capital
3 population
4 area
5 sea
```

Метод **.iterrows()** возвращает индекс строки и ее содержимое в формате Series.

```
1 # прервем цикл после первой итерации с помощью break
2 for index, row in countries.iterrows():
3     print(index)
4     print(row)
5     print('...')
6     print(type(row))
7     break
```

```
1 CN
2 country      China
3 capital      Beijing
4 population    1400
5 area          9.6
6 sea           1
7 Name: CN, dtype: object
8 ...
9 <class 'pandas.core.series.Series'>
```

Получить доступ к элементам одной строки можно по индексу объекта Series (то есть по

названиям столбцов исходного датафрейма).

```
1 for _, row in countries.iterrows():
2     # например, сформируем вот такое предложение
3     print(row['capital'] + ' is the capital of ' + row['country'])
4     break
```

```
1 Beijing is the capital of China
```

Доступ к столбцам

В отличие от объекта Series в датафрейме через квадратные скобки мы получаем доступ к его столбцам.

```
1 # выведем столбец capital датафрейма countries
2 countries['capital']
```

```
1 CN      Beijing
2 VN      Hanoi
3 GB      London
4 RU      Moscow
5 AR      Buenos Aires
6 BO      Sucre
7 ZA      Pretoria
8 Name: capital, dtype: object
```

Доступ к столбцу можно также получить, указав название датафрейма и, через точку, название искомого столбца.

```
1 # однако в этом случае название не должно содержать пробелов
2 countries.capital
```

```
1 CN      Beijing
2 VN      Hanoi
3 GB      London
4 RU      Moscow
5 AR      Buenos Aires
6 BO      Sucre
7 ZA      Pretoria
8 Name: capital, dtype: object
```

Как уже было сказано, отдельные столбцы в датафрейме имеют тип данных Series.

```
1 type(countries.capital)
```

```
1 pandas.core.series.Series
```

Для того чтобы получить доступ к столбцам и при этом на выходе сформировать датафрейм, необходимо использовать двойные скобки.

```
1 # логика здесь в том, что внутренние скобки - это список,
2 # внешние - оператор индексации
3 countries[['capital']]
```

capital	
CN	Beijing
VN	Hanoi
GB	London
RU	Moscow

AR	Buenos Aires
BO	Sucre
ZA	Pretoria

Так мы можем получить доступ к нескольким столбцам.

```
1 | countries[['capital', 'area']]
```

	capital	area
CN	Beijing	9.6
VN	Hanoi	0.3
GB	London	0.2
RU	Moscow	17.1
AR	Buenos Aires	2.8
BO	Sucre	1.1
ZA	Pretoria	1.2

Доступ к столбцам можно также получить с помощью **метода .filter()**, передав параметру *items* список с необходимыми нам столбцами.

```
1 | countries.filter(items = ['capital', 'population'])
```

	capital	population
CN	Beijing	1400
VN	Hanoi	97
GB	London	67
RU	Moscow	144
AR	Buenos Aires	45
BO	Sucre	12
ZA	Pretoria	59

Доступ к строкам

Получить доступ к строкам можно через срез индекса.

```
1 | # выведем строки со второй по пятую (не включительно)
2 | countries[1:5]
```

	country	capital	population	area	sea
VN	Vietnam	Hanoi	97	0.3	1
GB	United Kingdom	London	67	0.2	1
RU	Russia	Moscow	144	17.1	1
AR	Argentina	Buenos Aires	45	2.8	1

Методы .loc[] и .iloc[]

Метод .loc[]

Метод .loc[] позволяет получить доступ к строкам и столбцам через их названия (label-based location). Например, выведем первые три строки и первые три столбца датафрейма про страны.

```
1 # для этого передадим методу .loc[] два списка:  
2 # с индексами строк и названиями столбцов  
3 countries.loc[['CN', 'RU', 'VN'], ['capital', 'population', 'area']]
```

	capital	population	area
CN	Beijing	1400	9.6
RU	Moscow	144	17.1
VN	Hanoi	97	0.3

Через двоеточие, как и в Numpy, мы можем вывести все строки или все столбцы датафрейма.

```
1 # например, выведем все строки датафрейма  
2 countries.loc[:, ['capital', 'population', 'area']]
```

	capital	population	area
CN	Beijing	1400	9.6
VN	Hanoi	97	0.3
GB	London	67	0.2
RU	Moscow	144	17.1
AR	Buenos Aires	45	2.8
BO	Sucre	12	1.1
ZA	Pretoria	59	1.2

Метод .loc[] также поддерживает значения Boolean.

```
1 # например, выведем все строки и только последний столбец,  
2 # передав список соответствующих логических значений  
3 countries.loc[:, [False, False, False, False, True]]
```

	sea
CN	1
VN	1
GB	1
RU	1
AR	1
BO	0
ZA	1

Метод .get_loc() позволяет узнать порядковый номер (начиная с нуля) строки или столбца по

их индексу и названию соответственно. Для того чтобы узнать порядковый номер строки, метод `.get_loc()` нужно использовать совместно с атрибутом `index`.

```
1 # выведем номер строки с индексом RU
2 countries.index.get_loc('RU')
```

```
1 '3'
```

Теперь выведем номер столбца `country`. Для этого нам понадобится атрибут `columns`.

```
1 countries.columns.get_loc('country')
```

```
1 '0'
```

Метод `.iloc[]`

Метод `.iloc[]` действует примерно также, как и `.loc[]`, с тем отличием, что он основывается на числовом индексе (integer-based location).

```
1 # теперь в списки мы передаем номера строк и столбцов,
2 # нумерация начинается с нуля
3 countries.iloc[[0, 3, 5], [0, 1, 2]]
```

	country	capital	population
CN	China	Beijing	1400
RU	Russia	Moscow	144
BO	Bolivia	Sucre	12

В методе `.iloc[]` можно использовать срезы.

```
1 # выведем первые три строки и последние два столбца
2 countries.iloc[:3, -2:]
```

	area	sea
CN	9.6	1
VN	0.3	1
GB	0.2	1

К столбцам удобно обращаться по их названиям, к строкам — по порядковому номеру (числовому индексу). Для того чтобы это реализовать, мы можем объединить двойные квадратные скобки и метод `.iloc[]`.

```
1 # вначале передадим названия столбцов в двойных скобках,
2 # затем номера строк через метод .iloc[]
3 countries[['population', 'area']].iloc[[0, 3]]
```

	population	area
CN	1400	9.6
RU	144	17.1

Многоуровневый индекс и методы .loc[] и .iloc[]

Поговорим про применение **методов .loc[] и .iloc[]** к датафрейму с многоуровневым индексом.

```
1 # вновь создадим датафрейм с многоуровневым индексом по строкам и столбцам
2 countries.index = custom_multindex
3 countries.columns = custom_multicols
4
5 countries
```

		names		data		
		country	capital	population	area	sea
region	code					
Asia	CN	China	Beijing	1400	9.6	1
	VN	Vietnam	Hanoi	97	0.3	1
Europe	GB	United Kingdom	London	67	0.2	1
	RU	Russia	Moscow	144	17.1	1
S. America	AR	Argentina	Buenos Aires	45	2.8	1
	BO	Bolivia	Sucre	12	1.1	0
Africa	ZA	South Africa	Pretoria	59	1.2	1

Для доступа к первой строке передадим **методу .loc[]** соответствующий двойной индекс.

```
1 countries.loc['Asia', 'CN']

1 names country      China
2      capital    Beijing
3 data  population    1400
4      area         9.6
5      sea          1
6 Name: (Asia, CN), dtype: object
```

Мы также можем передать значения в форме кортежей как для строк, так и для столбцов.

```
1 # выведем первую строку и столбцы с числовыми данными
2 countries.loc[('Asia', 'CN'), [('data', 'population'), ('data', 'area'), ('data', 'sea')]]

1 data  population    1400.0
2      area         9.6
3      sea          1.0
4 Name: (Asia, CN), dtype: float64
```

Доступ к строкам можно получить, указав внутри кортежа название региона и список с кодами стран.

```
1 # например, выведем только азиатские страны
2 countries.loc[('Asia', ['CN', 'VN']), :]
```

		names		data		
		country	capital	population	area	sea
region	code					
Asia	CN	China	Beijing	1400	9.6	1
	VN	Vietnam	Hanoi	97	0.3	1

Внутри кортежа можно указать только регион. В этом случае мы получим все находящиеся в нем страны.

```
1 countries.loc(['Asia'], :)
```

names		data			
	country	capital	population	area	sea
code					
CN	China	Beijing	1400	9.6	1
VN	Vietnam	Hanoi	97	0.3	1

Аналогичным образом мы можем получить доступ к столбцам.

```
1 countries.loc[:, [('names', 'country'), ('data', 'population')]]
```

names		data	
	country	population	
region	code		
Asia	CN	China	1400
	VN	Vietnam	97
Europe	GB	United Kingdom	67
	RU	Russia	144
S. America	AR	Argentina	45
	BO	Bolivia	12
Africa	ZA	South Africa	59

Метод .iloc[], при этом, игнорирует структуру иерархического индекса и использует простой числовой индекс.

```
1 # получим доступ к четвертой строке и третьему, четвертому и пятому столбцам
2 countries.iloc[3, [2, 3, 4]]
```

```
1 data    population    144.0
2         area         17.1
3         sea          1.0
4 Name: (Europe, RU), dtype: float64
```

Метод .xs()

Метод .xs() (от англ. cross-section, срез) позволяет получить доступ к определенному уровню иерархического индекса. Начнем со строк.

```
1 # выберем Европу из уровня region
2 # axis = 0 указывает, что мы берем строки
3 countries.xs('Europe', level = 'region', axis = 0)
```

names		data			
	country	capital	population	area	sea

	country	capital	population	area	sea
code					
GB	United Kingdom	London	67	0.2	1
RU	Russia	Moscow	144	17.1	1

Перейдем к столбцам и выберем `names` и `country` (по сути мы выбираем столбец `country`) на первом и втором уровнях соответственно.

```
1 # levels указывает, на каких уровнях искать названия столбцов
2 # параметр axis = 1 говорит о том, что мы имеем дело со столбцами
3 countries.xs('names', 'country', level = [0, 1], axis = 1)
```

	names	
	country	
region	code	
Asia	CN	China
	VN	Vietnam
Europe	GB	United Kingdom
	RU	Russia
S. America	AR	Argentina
	BO	Bolivia
Africa	ZA	South Africa

Кроме того, мы можем соединить два метода `.xs()` и одновременно получить доступ и к строкам, и к столбцам.

```
1 # в данном случае мы можем не указывать level, потому что
2 # Europe и names находятся во внешних индексах, которые в level указаны по умолчанию
3 countries.xs('Europe', axis = 0).xs('names', axis = 1)
```

	country	capital
code		
GB	United Kingdom	London
RU	Russia	Moscow

Вернем датафрейму одноуровневый индекс.

```
1 # обновим атрибуты index и columns
2 countries.index = custom_index
3 countries.columns = custom_cols
4
5 # посмотрим на исходный датафрейм
6 countries
```

	country	capital	population	area	sea
CN	China	Beijing	1400	9.6	1
VN	Vietnam	Hanoi	97	0.3	1
GB	United Kingdom	London	67	0.2	1

RU	Russia	Moscow	144	17.1	1
AR	Argentina	Buenos Aires	45	2.8	1
BO	Bolivia	Sucre	12	1.1	0
ZA	South Africa	Pretoria	59	1.2	1

Метод .at[]

Метод `.at[]` подходит для извлечения или записи *одного* значения датафрейма.

```
1 | countries.at['CN', 'capital']
```

```
1 | 'Beijing'
```

Фильтры

Логическая маска

Как и в случае с массивом Numpy, мы можем создать логическую маску (Boolean mask) датафрейма, в которой нежелательные значения будут помечены как False.

```
1 | # создадим логическую маску для стран с населением больше миллиарда человек
2 | countries.population > 1000
```

```
1 | CN      True
2 | VN      False
3 | GB      False
4 | RU      False
5 | AR      False
6 | BO      False
7 | ZA      False
8 | Name: population, dtype: bool
```

Применив эту маску к исходному датафрейму, мы получим только те значения, которые помечены как True.

```
1 | # применим логическую маску к исходному датафрейму
2 | countries[countries.population > 1000]
```

	country	capital	population	area	sea
CN	China	Beijing	1400	9.6	1

Мы можем применить к датафрейму две маски, объединенные, например, логическим оператором И.

```
1 | # отфильтруем датафрейм по критериям численности населения и площади
2 | countries[(countries.population > 50) & (countries.area < 2)]
```

	country	capital	population	area	sea
VN	Vietnam	Hanoi	97	0.3	1
GB	United Kingdom	London	67	0.2	1
ZA	South Africa	Pretoria	59	1.2	1

Приведу еще один вариант синтаксиса. Его удобно использовать, если у вас есть множество условий, по которым вы хотите отфильтровать датафрейм.

```
1 # вначале создаем нужные нам маски
2 population_mask = countries.population > 70
3 area_mask = countries.population < 50
4
5 # затем объединяем их по необходимым условиям (в данном случае ИЛИ)
6 mask = population_mask | area_mask
7 # и применяем маску к исходному датафрейму
8 countries[mask]
```

	country	capital	population	area	sea
CN	China	Beijing	1400	9.6	1
VN	Vietnam	Hanoi	97	0.3	1
RU	Russia	Moscow	144	17.1	1
AR	Argentina	Buenos Aires	45	2.8	1
BO	Bolivia	Sucre	12	1.1	0

Метод .query()

Метод `.query()` позволяет задавать условие фильтрации «своими словами».

```
1 # например, выберем страны с населением более 50 млн. человек И
2 # площадью менее двух млн. кв. километров
3 countries.query('population > 50 and area < 2')
```

	country	capital	population	area	sea
VN	Vietnam	Hanoi	97	0.3	1
GB	United Kingdom	London	67	0.2	1
ZA	South Africa	Pretoria	59	1.2	1

При фильтрации по строковым значениям (которые записываются в кавычках), необходимо использовать различающиеся кавычки: двойные для метода `.query()` и одинарные для строкового значения или наоборот.

```
1 # выведем данные по Великобритании
2 countries.query("country == 'United Kingdom'")
```

	country	capital	population	area	sea
GB	United Kingdom	London	67	0.2	1

Другие способы фильтрации

С помощью метода `.isin()` мы можем проверить наличие нескольких значений в определенном столбце, а затем использовать результат в качестве логической маски.


```

1 # найдем строки, в которых в столбце capital присутствуют следующие значения
2 keyword_list = ['Beijing', 'Moscow', 'Hanoi']
3
4 countries[countries.capital.isin(keyword_list)]

```

	country	capital	population	area	sea
CN	China	Beijing	1400	9.6	1
VN	Vietnam	Hanoi	97	0.3	1
RU	Russia	Moscow	144	17.1	1

Похожим образом можно использовать метод `.startswith()`.

```

1 # например, для нахождения стран, НЕ начинающихся с буквы "А"
2 countries[~countries.country.str.startswith('A')]

```

	country	capital	population	area	sea
CN	China	Beijing	1400	9.6	1
VN	Vietnam	Hanoi	97	0.3	1
GB	United Kingdom	London	67	0.2	1
RU	Russia	Moscow	144	17.1	1
BO	Bolivia	Sucre	12	1.1	0
ZA	South Africa	Pretoria	59	1.2	1

Метод `.nlargest()` позволяет найти строки с наибольшим значением в определенном столбце.

```

1 # например, возьмем три строки с наибольшими значениями по столбцу population
2 countries.nlargest(3, 'population')

```

	country	capital	population	area	sea
CN	China	Beijing	1400	9.6	1
RU	Russia	Moscow	144	17.1	1
VN	Vietnam	Hanoi	97	0.3	1

Метод `.nsmallest()` аналогичным образом находит наименьшие значения.

Метод `.argmax()` выводит индекс строки, в которой в определенном столбце содержится наибольшее значение.

```

1 # например, предположим, что мы хотим найти индекс страны с наибольшей площадью
2 countries.area.argmax()

```

```

1 3

```

Посмотрим, какой стране соответствует этот индекс.

```

1 # напомним, что двойные скобки позволяют вывести DataFrame, а не Series
2 countries.iloc[[countries.area.argmax()]]

```

	country	capital	population	area	sea
--	---------	---------	------------	------	-----

RU	Russia	Moscow	144	17.1	1
----	--------	--------	-----	------	---

Метод .argmin() выводит, соответственно, индекс строки, в которой содержится наименьшее для заданного столбца значение.

Вспомним, что в метод **.loc[]** можно передать тип данных Boolean. Используем это свойство для создания фильтра.

```
1 # выведем страны с населением более 90 млн. человек
2 countries.loc[countries.population > 90, :]
```

	country	capital	population	area	sea
CN	China	Beijing	1400	9.6	1
VN	Vietnam	Hanoi	97	0.3	1
RU	Russia	Moscow	144	17.1	1

Метод .filter(), если использовать параметр `like`, ищет совпадения с искомой фразой в индексе (если `axis = 0`) или названии столбцов (если `axis = 1`).

```
1 # найдем строки, в которых в индексе есть буквосочетание "ZA"
2 countries.filter(like = 'ZA', axis = 0)
```

	country	capital	population	area	sea
ZA	South Africa	Pretoria	59	1.2	1

Очевидно, если указать параметр `axis = 1`, выведется индекс, потому что именно в нем есть буквосочетание ZA.

Сортировка

Для сортировки значений в датафрейме библиотеки Pandas используется **метод .sort_values()**. Этому методу мы можем передать, среди прочих, следующие параметры:

- **by** — по какому столбцу или столбцам вести сортировку;
- **inplace** — сохранять ли результат;
- **ascending** — в восходящем или нисходящем порядке сортировать.

Рассмотрим два примера.

```
1 # выполним сортировку по столбцу population, не сохраняя изменений,
2 # в возрастающем порядке (значение по умолчанию)
3 countries.sort_values(by = 'population', inplace = False, ascending = True)
```

	country	capital	population	area	sea
BO	Bolivia	Sucre	12	1.1	0
AR	Argentina	Buenos Aires	45	2.8	1
ZA	South Africa	Pretoria	59	1.2	1

GB	United Kingdom	London	67	0.2	1
VN	Vietnam	Hanoi	97	0.3	1
RU	Russia	Moscow	144	17.1	1
CN	China	Beijing	1400	9.6	1

```
1 # теперь отсортируем по двум столбцам в нисходящем порядке
2 countries.sort_values(by = ['area', 'population'], inplace = False, ascending = False)
```

	country	capital	population	area	sea
RU	Russia	Moscow	144	17.1	1
CN	China	Beijing	1400	9.6	1
AR	Argentina	Buenos Aires	45	2.8	1
ZA	South Africa	Pretoria	59	1.2	1
BO	Bolivia	Sucre	12	1.1	0
VN	Vietnam	Hanoi	97	0.3	1
GB	United Kingdom	London	67	0.2	1

Кроме того, можно отсортировать строки по индексу с помощью метода `.sort_index()`.

```
1 countries.sort_index()
```

	country	capital	population	area	sea
AR	Argentina	Buenos Aires	45	2.8	1
BO	Bolivia	Sucre	12	1.1	0
CN	China	Beijing	1400	9.6	1
GB	United Kingdom	London	67	0.2	1
RU	Russia	Moscow	144	17.1	1
VN	Vietnam	Hanoi	97	0.3	1
ZA	South Africa	Pretoria	59	1.2	1

На этом завершим сегодняшнее занятие.

Подведем итог

На первом занятии по анализу и обработке данных мы рассмотрели, из каких этапов состоит решение задачи машинного обучения и начали углубленное изучение библиотеки Pandas.

В частности, мы узнали про типы объектов библиотеки Pandas (датафрейм и Series), способы их создания, структуру и свойства. Кроме того, мы рассмотрели возможности доступа к строкам, столбцам и элементам датафрейма.