

Преобразование датафреймов

Материалы > Анализ и обработка данных

Продолжим работу по изучению библиотеки Pandas. Сегодня мы поговорим про возможности изменения и соединения датафреймов, а также про способы группировки данных.

Хотя сами по себе эти навыки являются скорее вспомогательными, овладеть ими совершенно необходимо, если вы хотите быстро и эффективно строить модели ML.

Откроем блокнот к этому занятию

Изменение датафрейма

Вернемся к тому датафрейму, с которым мы работали на прошлом занятии.

```
1 # создадим несколько списков и массивов Numpy с информацией о семи странах мира
2 country = np.array(['China', 'Vietnam', 'United Kingdom', 'Russia', 'Argentina', 'Boli
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)
7
8 # кроме того создадим список кодов стран, которые станут индексом датафрейма
9 custom_index = ['CN', 'VN', 'GB', 'RU', 'AR', 'BO', 'ZA']
10
11 # создадим пустой словарь
12 countries_dict = {}
13
14 # превратим эти списки в значения словаря,
15 # одновременно снабдив необходимыми ключами
16 countries_dict['country'] = country
17 countries_dict['capital'] = capital
18 countries_dict['population'] = population
19 countries_dict['area'] = area
20 countries_dict['sea'] = sea
21
22 # создадим датафрейм
23 countries = pd.DataFrame(countries_dict, index = custom_index)
24 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

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

Посмотрим, как мы можем преобразовать этот датафрейм.

Копирование датафрейма

Метод .copy()

В первую очередь поговорим про важную особенность при копировании датафрейма. Вначале создадим копию датафрейма с помощью простого присвоения этого объекта новой переменной.

```
1 countries_new = countries
```

Теперь удалим строку с данными про Аргентину, а после этого выведем исходный датафрейм.

```
1 countries_new.drop(labels = 'AR', axis = 0, inplace = True)
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
BO	Bolivia	Sucre	12	1.1	0
ZA	South Africa	Pretoria	59	1.2	1

Как вы видите, изменения коснулись и его. По этой причине для создания полноценной копии лучше использовать **метод .copy()**.

```
1 # в первую очередь вернем Аргентину в исходный датафрейм countries
2 countries = pd.DataFrame(countries_dict, index = custom_index)
3
4 # создадим копию, на этот раз с помощью метода .copy()
5 countries_new = countries.copy()
6
7 # вновь удалим запись про Аргентину
8 countries_new.drop(labels = 'AR', axis = 0, inplace = True)
9
10 # выведем исходный датафрейм
11 countries
```

	country	capital	population	area	sea
CN	China	Beijing	1400	9.6	1
VN	Vietnam	Hanoi	97	0.3	1
GB	United Kinadom	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

Такой тип копирования (без изменения исходного датафрейма) может быть полезен, например, когда мы хотим взять датафрейм, протестировать какие-то гипотезы, построить одну или несколько моделей, но при этом не трогать исходные данные.

Про параметр inplace

Создадим несложный датафрейм из вложенных списков.

```
1 df = pd.DataFrame([[1, 1, 1],
2                   [2, 2, 2],
3                   [3, 3, 3]],
4                   columns = ['A', 'B', 'C'])
5
6 df
```

	A	B	C
0	1	1	1
1	2	2	2
2	3	3	3

Как понять, сохраняется ли изменение после применение определенного метода или нет? Если метод выдает датафрейм, изменение не сохраняется.

```
1 # попробуем удалить столбец A
2 df.drop(labels = ['A'], axis = 1)
```

	B	C
0	1	1
1	2	2
2	3	3

```
1 # проверим, сохранилось ли изменение
2 df
```

	A	B	C
0	1	1	1
1	2	2	2
2	3	3	3

При этом если метод выдает None, изменение постоянно.

```
1 # изменим параметр inplace на True
2 print(df.drop(labels = ['A'], axis = 1, inplace = True))
```

```
1 None
```

```
1 # проверим
2 df
```

	B	C
0	1	1
1	2	2
2	3	3

По этой причине нельзя использовать `inplace = True` и записывать результат в переменную одновременно.

```
1 df = df.drop(labels = ['B'], axis = 1, inplace = True)
2 print(df)
```

```
1 None
```

В этом случае мы записываем `None` в переменную `df`.

Столбцы датафрейма

Именованние столбцов при создании датафрейма

Создадим список с названиями столбцов на кириллице и транспонированный массив Numpy с данными о странах.

```
1 custom_columns = ['страна', 'столица', 'население', 'площадь', 'море']
2 arr = np.array([country, capital, population, area, sea]).T
3 arr
```

```
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='<U32')
```

После этого создадим датафрейм с помощью функции `pd.DataFrame()` с параметром `columns`, в который передадим названия столбцов.

```
1 countries = pd.DataFrame(data = arr,
2                           index = custom_index,
3                           columns = custom_columns)
4
5 countries
```

	1	2	2	2	страна	столица	население	площадь	море
CN	3	3	3	3	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

Вернем прежние названия столбцов.

```
1 | countries.columns = ['country', 'capital', 'population', 'area', 'sea']
```

Переименование столбцов

Для того чтобы переименовать отдельные столбцы, можно воспользоваться **методом .rename()**. Внутри этого метода в параметр `columns` мы передаем словарь, где ключами будут текущие названия столбцов, а значениями соответствующие им новые названия.

```
1 | # переименуем столбец capital на city
2 | countries.rename(columns = {'capital': 'city'}, inplace = True)
3 | countries
```

	country	city	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

Тип данных в столбце

Все значения одного столбца датафрейма всегда имеют один и тот же тип данных. Посмотреть тип данных каждого из столбцов можно с помощью **атрибута .dtypes**.

```
1 | countries.dtypes
```

```
1 | country      object
2 | city         object
3 | population   object
4 | area         object
5 | sea          object
6 | dtype: object
```

Изменение типа данных

Преобразовать тип данных столбца можно с помощью **метода .astype()**. Этот метод можно применить к конкретному столбцу.

```
1 | # преобразуем тип данных столбца population в int
2 | countries.population = countries.population.astype('int')
```

Кроме того, мы можем преобразовать тип данных сразу в нескольких столбцах. Для этого применим метод `.astype()` ко всему датафрейму. Самому методу мы передадим словарь, где ключами будут названия столбцов, а значениями — соответствующий им желаемый тип данных.

```
1 # изменим тип данных в столбцах area и sea
2 countries = countries.astype({'area': 'float', 'sea' : 'category'})
```

Посмотрим на результат.

```
1 countries.dtypes

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

Тип данных category

Обратите внимание на новый для нас тип данных `category`. Во многом он похож на факторную переменную в R.

```
1 # в category содержится информация об имеющихся в столбце категориях
2 countries.sea
```

```
1 CN      1
2 VN      1
3 GB      1
4 RU      1
5 AR      1
6 BO      0
7 ZA      1
8 Name: sea, dtype: category
9 Categories (2, object): ['0', '1']
```

Тип `category` мы рассмотрим более подробно на занятии по кодированию категориальных переменных.

Помимо упомянутых типов данных, нам также знаком объект `datetime`, который используется для работы с временными рядами. Мы снова обратимся к нему на занятии по очистке данных.

Фильтр столбцов по типу данных

Выбрать столбцы в зависимости от типа содержащихся в них данных можно с помощью метода `.select_dtypes()`. Включить определенные типы данных можно с помощью параметра `include`.

```
1 # выберем только типы данных int и float
2 countries.select_dtypes(include = ['int64', 'float64'])
```

	population	area
CN	1400	9.6
VN	97	0.3

GB	67	0.2
RU	144	17.1
AR	45	2.8
BO	12	1.1
ZA	59	1.2

Исключить определенные типы данных можно через *exclude*.

```
1 # выберем все типы данных, кроме object и category
2 countries.select_dtypes(exclude = ['object', 'category'])
```

	population	area
CN	1400	9.6
VN	97	0.3
GB	67	0.2
RU	144	17.1
AR	45	2.8
BO	12	1.1
ZA	59	1.2

Добавление строк и столбцов

Добавление строк

Метод `._append()` + словарь

Для добавления строк в первую очередь используется **метод `.append()`**. С его помощью строку можно добавить из питоновского словаря.

```
1 # создадим словарь с данными Канады и добавим его в датафрейм
2 dict_ = {'country': 'Canada', 'city': 'Ottawa', 'population': 38, 'area': 10, 'sea' : '1'
3
4 # словарь можно добавлять только если ignore_index = True
5 countries = countries._append(dict_, ignore_index = True)
6 countries
```

	country	city	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
7	Canada	Ottawa	38	10.0	1

Метод ._append() + Series

Мы также можем добавить строки в виде объекта Series.

```
1 # причем, если передать список из Series, можно добавить сразу несколько строк
2 list_of_series = [pd.Series(['Spain', 'Madrid', 47, 0.5, 1], index = countries.columns),
3                  pd.Series(['Netherlands', 'Amsterdam', 17, 0.04, 1], index = countries.columns)]
4
5 # нам по-прежнему необходим параметр ignore_index = True
6 countries._append(list_of_series, ignore_index = True)
7 countries
```

	country	city	population	area	sea
0	China	Beijing	1400	9.60	1
1	Vietnam	Hanoi	97	0.30	1
2	United Kingdom	London	67	0.20	1
3	Russia	Moscow	144	17.10	1
4	Argentina	Buenos Aires	45	2.80	1
5	Bolivia	Sucre	12	1.10	0
6	South Africa	Pretoria	59	1.20	1
7	Canada	Ottawa	38	10.00	1
8	Spain	Madrid	47	0.50	1
9	Netherlands	Amsterdam	17	0.04	1

Метод ._append() + другой датафрейм

Новая строка может также содержаться в другом датафрейме.

```
1 # новая строка может также содержаться в другом датафрейме
2 # обратите внимание, что числовые значения мы помещаем в списки
3 peru = pd.DataFrame({'country' : 'Peru',
4                      'city' : 'Lima',
5                      'population': [33],
6                      'area' : [1.3],
7                      'sea' : [1]})
8 peru
```

	country	city	population	area	sea
0	Peru	Lima	33	1.3	1

```
1 # перед добавлением выберем первую строку с помощью метода .iloc[]
2 countries._append(peru.iloc[0], ignore_index = True)
```

	country	city	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

4	Argentina	Buenos Aires	45	2.8	1
5	Bolivia	Sucre	12	1.1	0
6	South Africa	Pretoria	59	1.2	1
7	Canada	Ottawa	38	10.0	1
8	Peru	Lima	33	1.3	1

Использование .iloc[]

Если вновь вывести наш датафрейм `countries`, мы увидим, что данные об Испании, Нидерландах и Перу не сохранились.

```
1 # для этого нам надо было либо перезаписать результат метода ._append() в переменную с
2 # либо использовать параметр inplace = True.
3 countries
```

	country	city	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
7	Canada	Ottawa	38	10.0	1

Добавим данные об этих странах на постоянной основе с помощью **метода .iloc[]** и посмотрим на результат.

```
1 countries.iloc[[5, 6, 7]] = [['Spain', 'Madrid', 47, 0.5, 1],
2                               ['Netherlands', 'Amsterdam', 17, 0.04, 1],
3                               ['Peru', 'Lima', 33, 1.3, 1]]
4
5 countries
```

	country	city	population	area	sea
0	China	Beijing	1400	9.60	1
1	Vietnam	Hanoi	97	0.30	1
2	United Kingdom	London	67	0.20	1
3	Russia	Moscow	144	17.10	1
4	Argentina	Buenos Aires	45	2.80	1
5	Spain	Madrid	47	0.50	1
6	Netherlands	Amsterdam	17	0.04	1
7	Peru	Lima	33	1.30	1

Обратите внимание, что строки добавились строго на те индексы, которые были указаны в

методе `.iloc[]` (т.е. 5, 6 и 7) и заменили данные, ранее находившиеся на этих индексах (Боливия, Южная Африка и Канада).

В версии Pandas 2.0.0, которая была опубликована 3 апреля 2023 года, метод `.append()` был удален, и применение нижнего подчеркивания, хотя и позволяет выполнить присоединение строк, не является удачным решением.

Разработчики рекомендуют использовать **метод `.concat()`**. О нем мы поговорим ниже.

Добавление столбцов

Объявление нового столбца

Новый столбец датафрейма можно просто объявить и сразу добавить в него необходимые данные.

```
1 # например, добавим данные о плотности населения
2 countries['pop_density'] = [153, 49, 281, 9, 17, 94, 508, 26]
3 countries
```

	country	city	population	area	sea	pop_density
0	China	Beijing	1400	9.60	1	153
1	Vietnam	Hanoi	97	0.30	1	49
2	United Kingdom	London	67	0.20	1	281
3	Russia	Moscow	144	17.10	1	9
4	Argentina	Buenos Aires	45	2.80	1	17
5	Spain	Madrid	47	0.50	1	94
6	Netherlands	Amsterdam	17	0.04	1	508
7	Peru	Lima	33	1.30	1	26

Метод `.insert()`

Добавим столбец с кодами стран с помощью **метода `.insert()`**.

```
1 countries.insert(loc = 1, # это будет второй по счету столбец
2                  column = 'code', # название столбца
3                  value = ['CN', 'VN', 'GB', 'RU', 'AR', 'ES', 'NL', 'PE']) # значения с
```

Обратите внимание, метод `.insert()` по умолчанию (без перезаписи в переменную или параметра `inplace = True`) сохраняет результат.

```
1 # посмотрим на результат
2 countries
```

	country	code	city	population	area	sea	pop_density
0	China	CN	Beijing	1400	9.60	1	153
1	Vietnam	VN	Hanoi	97	0.30	1	49
2	United Kingdom	GB	London	67	0.20	1	281
3	Russia	RU	Moscow	144	17.10	1	9
4	Argentina	AR	Buenos Aires	45	2.80	1	17
5	Spain	ES	Madrid	47	0.50	1	94
6	Netherlands	NL	Amsterdam	17	0.04	1	508
7	Peru	PE	Lima	33	1.30	1	26

Теперь рассмотрим несколько способов добавления столбца с рассчитанным значением.

Метод .assign()

Создадим столбец `area_miles`, в который поместим площадь в милях. Вначале используем метод `.assign()`.

```
1 countries = countries.assign(area_miles = countries.area / 2.59).round(2)
2 countries
```

	country	code	city	population	area	sea	pop_density	area_miles
0	China	CN	Beijing	1400	9.60	1	153	3.71
1	Vietnam	VN	Hanoi	97	0.30	1	49	0.12
2	United Kingdom	GB	London	67	0.20	1	281	0.08
3	Russia	RU	Moscow	144	17.10	1	9	6.60
4	Argentina	AR	Buenos Aires	45	2.80	1	17	1.08
5	Spain	ES	Madrid	47	0.50	1	94	0.19
6	Netherlands	NL	Amsterdam	17	0.04	1	508	0.02
7	Peru	PE	Lima	33	1.30	1	26	0.50

Удалим этот столбец, чтобы рассмотреть другие методы.

```
1 countries.drop(labels = 'area_miles', axis = 1, inplace = True)
```

Можно сложнее

Мы можем усложнить код и добиться такого же результата, применив методы `.iterrows()` и `.iloc[]`.

```
1 # выведем индекс и содержание строк
2 for index, row in countries.iterrows():
3     # запишем для каждой строки (index) в новый столбец area_miles
4     # округленное значение площади row.area в милях
5     countries.loc[index, 'area_miles'] = np.round(row.area / 2.59, 2)
6
7 # посмотрим на результат
8 countries
```

	country	code	city	population	area	sea	pop_density	area_miles
0	China	CN	Beijing	1400	9.60	1	153	3.71
1	Vietnam	VN	Hanoi	97	0.30	1	49	0.12
2	United Kingdom	GB	London	67	0.20	1	281	0.08
3	Russia	RU	Moscow	144	17.10	1	9	6.60
4	Argentina	AR	Buenos Aires	45	2.80	1	17	1.08
5	Spain	ES	Madrid	47	0.50	1	94	0.19
6	Netherlands	NL	Amsterdam	17	0.04	1	508	0.02
7	Peru	PE	Lima	33	1.30	1	26	0.50

```
1 # снова удалим этот столбец
2 countries.drop(labels = 'area_miles', axis = 1, inplace = True)
```

Можно проще

При этом конечно есть гораздо более простой способ добавления нового столбца.

```
1 # мы можем объявить столбец и присвоить ему нужно нам значение
2 countries['area_miles'] = (countries.area / 2.59).round(2)
3 countries
```

	country	code	city	population	area	sea	pop_density	area_miles
0	China	CN	Beijing	1400	9.60	1	153	3.71
1	Vietnam	VN	Hanoi	97	0.30	1	49	0.12
2	United Kingdom	GB	London	67	0.20	1	281	0.08
3	Russia	RU	Moscow	144	17.10	1	9	6.60
4	Argentina	AR	Buenos Aires	45	2.80	1	17	1.08
5	Spain	ES	Madrid	47	0.50	1	94	0.19
6	Netherlands	NL	Amsterdam	17	0.04	1	508	0.02
7	Peru	PE	Lima	33	1.30	1	26	0.50

Удаление строк и столбцов

Удаление строк

Для удаления строк можно использовать **метод .drop()** с параметрами *labels* (индекс удаляемых строк) и *axis = 0*.

```
1 # удалим строки с индексом 0 и 1
2 countries.drop(labels = [0, 1], axis = 0)
```

	country	code	city	population	area	sea	pop_density	area_miles
2	United Kingdom	GB	London	67	0.20	1	281	0.08
3	Russia	RU	Moscow	144	17.10	1	9	6.60
4	Argentina	AR	Buenos Aires	45	2.80	1	17	1.08
5	Spain	ES	Madrid	47	0.50	1	94	0.19
6	Netherlands	NL	Amsterdam	17	0.04	1	508	0.02
7	Peru	PE	Lima	33	1.30	1	26	0.50

5	Spain	ES	мадрид	47	0.50	1	94	0.19
6	Netherlands	NL	Amsterdam	17	0.04	1	508	0.02
7	Peru	PE	Lima	33	1.30	1	26	0.50

Кроме того, можно использовать метод `.drop()` с единственным параметром `index`.

```
1 # удалим строки с индексом 5 и 7
2 countries.drop(index = [5, 7])
```

	country	code	city	population	area	sea	pop_density	area_miles
0	China	CN	Beijing	1400	9.60	1	153	3.71
1	Vietnam	VN	Hanoi	97	0.30	1	49	0.12
2	United Kingdom	GB	London	67	0.20	1	281	0.08
3	Russia	RU	Moscow	144	17.10	1	9	6.60
4	Argentina	AR	Buenos Aires	45	2.80	1	17	1.08
6	Netherlands	NL	Amsterdam	17	0.04	1	508	0.02

Мы также можем в параметр `index` передать индекс датафрейма через атрибут `index`.

```
1 # удалим четвертую строку
2 countries.drop(index = countries.index[4])
```

	country	code	city	population	area	sea	pop_density	area_miles
0	China	CN	Beijing	1400	9.60	1	153	3.71
1	Vietnam	VN	Hanoi	97	0.30	1	49	0.12
2	United Kingdom	GB	London	67	0.20	1	281	0.08
3	Russia	RU	Moscow	144	17.10	1	9	6.60
5	Spain	ES	Madrid	47	0.50	1	94	0.19
6	Netherlands	NL	Amsterdam	17	0.04	1	508	0.02
7	Peru	PE	Lima	33	1.30	1	26	0.50

С атрибутом датафрейма `index` мы можем делать срезы.

```
1 # удалим каждую вторую строку, начиная с четвертой с конца
2 countries.drop(index = countries.index[-4::2])
```

	country	code	city	population	area	sea	pop_density	area_miles
0	China	CN	Beijing	1400	9.6	1	153	3.71
1	Vietnam	VN	Hanoi	97	0.3	1	49	0.12
2	United Kingdom	GB	London	67	0.2	1	281	0.08
3	Russia	RU	Moscow	144	17.1	1	9	6.60
5	Spain	ES	Madrid	47	0.5	1	94	0.19
7	Peru	PE	Lima	33	1.3	1	26	0.50

Удаление столбцов

Параметры *labels* (номера столбцов) и *axis = 1* метода **.drop()** позволяют удалять столбцы.

```
1 # удалим столбцы area_miles и code
2 countries.drop(labels = ['area_miles', 'code'], axis = 1)
```

	country	city	population	area	sea	pop_density
0	China	Beijing	1400	9.60	1	153
1	Vietnam	Hanoi	97	0.30	1	49
2	United Kingdom	London	67	0.20	1	281
3	Russia	Moscow	144	17.10	1	9
4	Argentina	Buenos Aires	45	2.80	1	17
5	Spain	Madrid	47	0.50	1	94
6	Netherlands	Amsterdam	17	0.04	1	508
7	Peru	Lima	33	1.30	1	26

Такой же результат можно получить, передав список удаляемых столбцов в *параметр columns*.

```
1 # снова удалим столбцы area_miles и code
2 countries.drop(columns = ['area_miles', 'code'])
```

	country	city	population	area	sea	pop_density
0	China	Beijing	1400	9.60	1	153
1	Vietnam	Hanoi	97	0.30	1	49
2	United Kingdom	London	67	0.20	1	281
3	Russia	Moscow	144	17.10	1	9
4	Argentina	Buenos Aires	45	2.80	1	17
5	Spain	Madrid	47	0.50	1	94
6	Netherlands	Amsterdam	17	0.04	1	508
7	Peru	Lima	33	1.30	1	26

Через атрибут датафрейма **columns** мы можем передать номера удаляемых столбцов.

```
1 # удалим последний столбец (area_miles)
2 countries.drop(columns = countries.columns[-1])
```

	country	code	city	population	area	sea	pop_density
0	China	CN	Beijing	1400	9.60	1	153
1	Vietnam	VN	Hanoi	97	0.30	1	49
2	United Kingdom	GB	London	67	0.20	1	281
3	Russia	RU	Moscow	144	17.10	1	9
4	Argentina	AR	Buenos Aires	45	2.80	1	17
-	-	-	-	-	-	-	-

5	Spain	ES	Madrid	47	0.50	1	94
6	Netherlands	NL	Amsterdam	17	0.04	1	508
7	Peru	PE	Lima	33	1.30	1	26

Наконец удалим пятую строку и несколько столбцов разобранными выше способами и сохраним изменения.

```
1 countries.drop(index = 4, inplace = True)
2 countries.drop(columns = ['code', 'pop_density', 'area_miles'], inplace = True)
3 countries
```

	country	city	population	area	sea
0	China	Beijing	1400	9.60	1
1	Vietnam	Hanoi	97	0.30	1
2	United Kingdom	London	67	0.20	1
3	Russia	Moscow	144	17.10	1
5	Spain	Madrid	47	0.50	1
6	Netherlands	Amsterdam	17	0.04	1
7	Peru	Lima	33	1.30	1

Удаление по многоуровнему индексу

Давайте посмотрим, как удалять строки и столбцы в датафрейме с многоуровневым (иерархическим) индексом. Вначале вновь создадим соответствующий датафрейм.

```
1 # подготовим данные для многоуровневого индекса строк
2 rows = [('Asia', 'CN'),
3         ('Asia', 'VN'),
4         ('Europe', 'GB'),
5         ('Europe', 'RU'),
6         ('Europe', 'ES'),
7         ('Europe', 'NL'),
8         ('S. America', 'PE')]
9
10 # и столбцов
11 cols = [('names', 'country'),
12         ('names', 'city'),
13         ('data', 'population'),
14         ('data', 'area'),
15         ('data', 'sea')]
16
17 # создадим многоуровневый (иерархический) индекс
18 # для индекса строк добавим названия столбцов индекса через параметр names
19 custom_multindex = pd.MultiIndex.from_tuples(rows, names = ['region', 'code'])
20 custom_multicols = pd.MultiIndex.from_tuples(cols)
21
22 # поместим индексы в атрибуты index и columns датафрейма
23 countries.index = custom_multindex
24 countries.columns = custom_multicols
25
26 # посмотрим на результат
27 countries
```

names

data

region	country		city	population	area	sea
	code					
Asia	CN	China	Beijing	1400	9.60	1
	VN	Vietnam	Hanoi	97	0.30	1
Europe	GB	United Kingdom	London	67	0.20	1
	RU	Russia	Moscow	144	17.10	1
	ES	Spain	Madrid	47	0.50	1
	NL	Netherlands	Amsterdam	17	0.04	1
S. America	PE	Peru	Lima	33	1.30	1

Удаление строк

Вначале обратимся к строкам и удалим азиатский регион. Для этого воспользуемся методом `.drop()`, которому передадим параметр `labels = 'Asia'`. Кроме того, укажем, что удаляем именно строки (`axis = 0`) и что Азия находится в индексе под названием `region` (т.е. `level = 0`).

```
1 | countries.drop(labels = 'Asia', axis = 0, level = 0)
```

region	names		city	population	area	sea
	code	country				
Europe	GB	United Kingdom	London	67	0.20	1
	RU	Russia	Moscow	144	17.10	1
	ES	Spain	Madrid	47	0.50	1
	NL	Netherlands	Amsterdam	17	0.04	1
S. America	PE	Peru	Lima	33	1.30	1

Кроме того, строки можно удалять с помощью параметра `index` и указанием через `level`, по какому столбцу индекса мы будем искать удаляемую строку.

```
1 | # удалим запись о России по ее индексу в столбце code (т.е. level = 1)
2 | countries.drop(index = 'RU', level = 1)
```

region	names		city	population	area	sea
	code	country				
Asia	CN	China	Beijing	1400	9.60	1
	VN	Vietnam	Hanoi	97	0.30	1
Europe	GB	United Kingdom	London	67	0.20	1
	ES	Spain	Madrid	47	0.50	1
	NL	Netherlands	Amsterdam	17	0.04	1
S. America	PE	Peru	Lima	33	1.30	1

Удаление столбцов

Удаление столбцов датафрейма с многоуровневым индексом происходит аналогично строкам. Передадим методу `.drop()` параметры `labels`, `level` и `axis = 1` для удаления столбца по его наименованию (`labels`) на нужном нам уровне (`level`) индекса.

```
1 # удалим все столбцы в разделе names на нулевом уровне индекса столбцов
2 countries.drop(labels = 'names', level = 0, axis = 1)
```

data				
		population	area	sea
region	code			
Asia	CN	1400	9.60	1
	VN	97	0.30	1
Europe	GB	67	0.20	1
	RU	144	17.10	1
	ES	47	0.50	1
	NL	17	0.04	1
S. America	PE	33	1.30	1

Для удаления столбцов можно использовать параметр `columns` с указанием соответствующего уровня индекса (`level`) столбцов.

```
1 # например, удалим столбцы city и area на втором уровне (level = 1) индекса
2 countries.drop(columns = ['city', 'area'], level = 1)
```

names		data		
		country	population	sea
region	code			
Asia	CN	China	1400	1
	VN	Vietnam	97	1
Europe	GB	United Kingdom	67	1
	RU	Russia	144	1
	ES	Spain	47	1
	NL	Netherlands	17	1
S. America	PE	Peru	33	1

Обратите внимание, что удалению столбцов не помешал тот факт, что они находятся в разных индексах первого (`level = 0`) уровня, а именно `city` находится в `names`, а `population` — в `data`.

Применение функций

Библиотека Pandas позволяет использовать функции для изменения данных в датафрейме.

```
1 # создадим новый датафрейм с данными нескольких человек
2 people = pd.DataFrame({'name'      : ['Алексей', 'Иван', 'Анна', 'Ольга', 'Николай'],
3                        'gender'    : [1, 1, 0, 2, 1],
```

```
4         'age'       : [35, 20, 13, 28, 16],
5         'height'    : [180.46, 182.26, 165.12, 168.04, 178.68],
6         'weight'    : [73.61, 75.34, 50.22, 52.14, 69.72]
7     })
8
9     people
```

	name	gender	age	height	weight
0	Алексей	1	35	180.46	73.61
1	Иван	1	20	182.26	75.34
2	Анна	0	13	165.12	50.22
3	Ольга	2	28	168.04	52.14
4	Николай	1	16	178.68	69.72

Метод .map()

Предположим, что мы хотим явно прописать мужской (male) и женский (female) пол в наших данных. В этом случае можно воспользоваться **методом .map()**. Вначале создадим карту (map) того, как преобразовать существующие значения в новые.

```
1 # такая карта представляет собой питоновский словарь,
2 # где ключи - это старые данные, а значения - новые
3 gender_map = {0: 'female', 1: 'male'}
```

Применим эту карту к нужному нам столбцу датафрейма.

```
1 people['gender'] = people['gender'].map(gender_map)
2 people
```

	name	gender	age	height	weight
0	Алексей	male	35	180.46	73.61
1	Иван	male	20	182.26	75.34
2	Анна	female	13	165.12	50.22
3	Ольга	NaN	28	168.04	52.14
4	Николай	male	16	178.68	69.72

Те значения, которые в карте не указаны (в четвертой строке была ошибка, пол был помечен цифрой два), превращаются в NaN (not a number), пропущенное значение.

Обратите внимание, что словарь, в данном случае, по сути является функцией, которую мы применяем к значениям определенного столбца.

Кроме того замечу, что в большинстве случаев мы будем проводить обратное преобразование, превращая строковые категориальные значения в числовые.

В метод **.map()** мы можем передать и lambda-функцию.

```
1 # например, для того, чтобы выявить совершеннолетних и несовершеннолетних людей
2 people['age_group'] = people['age'].map(lambda x: 'adult' if x >= 18 else 'minor')
```

```
3 | people
```

	name	gender	age	height	weight	age_group
0	Алексей	male	35	180.46	73.61	adult
1	Иван	male	20	182.26	75.34	adult
2	Анна	female	13	165.12	50.22	minor
3	Ольга	NaN	28	168.04	52.14	adult
4	Николай	male	16	178.68	69.72	minor

Удалим только что созданный столбец age_group.

```
1 | people.drop(labels = 'age_group', axis = 1, inplace = True)
```

Вместо lambda-функции, например, для более сложных преобразований, можно использовать обычную собственную функцию.

```
1 | # обратите внимание, такая функция не допускает дополнительных параметров,  
2 | # только те данные, которые нужно преобразовать (age)  
3 | def get_age_group(age):  
4 |  
5 |     # например, мы не можем сделать threshold произвольным параметром  
6 |     threshold = 18  
7 |  
8 |     if age >= threshold:  
9 |         age_group = 'adult'  
10 |  
11 |     else:  
12 |         age_group = 'minor'  
13 |  
14 |     return age_group
```

Применим эту функцию к столбцу age.

```
1 | people['age_group'] = people['age'].map(get_age_group)  
2 | people
```

	name	gender	age	height	weight	age_group
0	Алексей	male	35	180.46	73.61	adult
1	Иван	male	20	182.26	75.34	adult
2	Анна	female	13	165.12	50.22	minor
3	Ольга	NaN	28	168.04	52.14	adult
4	Николай	male	16	178.68	69.72	minor

```
1 | # снова удалим созданный столбец  
2 | people.drop(labels = 'age_group', axis = 1, inplace = True)
```

Функция np.where()

Точно такой же результат мы можем получить, применив функцию **np.where()** библиотеки Numpy к нужному нам столбцу.

```
1 | # внутри функции np.where() три параметра: (1) условие,
```

```
2 | # (2) значение, если условие выдает True, (3) и значение, если условие выдает False
3 | people['age_group'] = np.where(people['age'] >= 18, 'adult', 'minor')
4 | people
```

	name	gender	age	height	weight	age_group
0	Алексей	male	35	180.46	73.61	adult
1	Иван	male	20	182.26	75.34	adult
2	Анна	female	13	165.12	50.22	minor
3	Ольга	NaN	28	168.04	52.14	adult
4	Николай	male	16	178.68	69.72	minor

Замечу, что такой способ превращения количественных данных в категориальные называется **binning** или **bucketing**. О нем мы подробно поговорим на занятии по преобразованию данных.

```
1 | # удалим созданный столбец
2 | people.drop(labels = 'age_group', axis = 1, inplace = True)
```

Метод .where()

Метод **.where()** библиотеки Pandas действует немного иначе. В нем мы прописываем условие, которое хотим применить к отдельному столбцу или всему датафрейму:

- если условие выполняется (т.е. оценивается как True), мы сохраняем текущее значение датафрейма;
- если условие не выполняется (False), то значение заменяется на новое, указанное в методе **.where()**.

Рассмотрим применение этого метода на примерах.

```
1 | # заменим возраст тех, кому меньше 18, на NaN
2 | people.age.where(people.age >= 18, other = np.nan)
```

```
1 | 0    35.0
2 | 1    20.0
3 | 2     NaN
4 | 3    28.0
5 | 4     NaN
6 | Name: age, dtype: float64
```

В примере выше возраст тех, кто не моложе 18 (True), остался без изменений. Для остальных (False) значение изменилось на пропущенное (параметр *other*). Обратите внимание, что тип данных этого столбца превратился во float. Это связано с тем, что в столбце появились пропущенные значения.

```
1 | # создадим матрицу из вложенных списков
2 | nums_matrix = [[-13, 7, 1],
3 |                [4, -2, 25],
4 |                [45, -3, 8]]
5 |
6 | # преобразуем в датафрейм
7 | # (матрица не обязательно должна быть массивом Numpy (!))
8 | nums = pd.DataFrame(nums_matrix)
9 | nums
```

	0	1	2
0	-13	7	1
1	4	-2	25
2	45	-3	8

```
1 # если число положительное (nums < 0 == True), оставим его без изменений
2 # если отрицательное (False), заменим на обратное (т.е. сделаем положительным)
3 nums.where(nums > 0, other = -nums)
```

	0	1	2
0	13	7	1
1	4	2	25
2	45	3	8

Здесь можно отметить два интересных момента:

- мы применили метод **.where()** ко всему датафрейму;
- в качестве значения, на которое нужно заменить текущее при False, мы использовали сами значения датафрейма, но со знаком минус (`-nums`).

Перейдем к методу **.apply()**.

Метод .apply()

Применение функции с аргументами

В отличие от **.map()**, метод **.apply()** позволяет передавать именованные аргументы в применяемую функцию.

```
1 # объявим функцию, которой можно передать не только значение возраста, но и порог,
2 # при котором мы будем считать человека совершеннолетним
3 def get_age_group(age, threshold):
4
5     if age >= int(threshold):
6         age_group = 'adult'
7     else:
8         age_group = 'minor'
9
10    return age_group
```

```
1 # применим эту функцию к столбцу age, выбрав в качестве порогового значения 21 год
2 people['age_group'] = people['age'].apply(get_age_group, threshold = 21)
3
4 # посмотрим на результат
5 people
```

	name	gender	age	height	weight	age_group
0	Алексей	male	35	180.46	73.61	adult
1	Иван	male	20	182.26	75.34	minor
2	Анна	female	13	165.12	50.22	minor
3	Ольга	NaN	28	168.04	52.14	adult

4	Николай	male	16	178.68	69.72	minor
---	---------	------	----	--------	-------	-------

Применение к столбцам

В метод `.apply()` можно передать уже имеющуюся в Питоне функцию, например, из библиотеки Numpy.

```
1 # заменим значения в столбцах height и weight на медиану по столбцам
2 people[['height', 'weight']] = people[['height', 'weight']].apply(np.median,
3                                                                    axis = 0)
4 people
```

	name	gender	age	height	weight	age_group
0	Алексей	male	35	178.68	69.72	adult
1	Иван	male	20	178.68	69.72	minor
2	Анна	female	13	178.68	69.72	minor
3	Ольга	NaN	28	178.68	69.72	adult
4	Николай	male	16	178.68	69.72	minor

Применение к строкам

Теперь поработаем со строками. Создадим функцию, которая считает индекс массы тела (body mass index, BMI) на основе веса и роста человека.

```
1 # внутри функции разделим вес на квадрат роста
2 def get_bmi(x):
3     bmi = x['weight'] / (x['height'] / 100) ** 2
4     return bmi
```

Теперь применим эту функцию к каждой строке и сохраним результат в новом столбце.

```
1 # для применения функции к строке используется параметр axis = 1
2 people['bmi'] = people.apply(get_bmi, axis = 1).round(2)
3 people
```

	name	gender	age	height	weight	age_group	bmi
0	Алексей	male	35	180.0	74.0	adult	22.84
1	Иван	male	20	182.0	75.0	minor	22.64
2	Анна	female	13	165.0	50.0	minor	18.37
3	Ольга	NaN	28	168.0	52.0	adult	18.42
4	Николай	male	16	179.0	70.0	minor	21.85

Метод .applymap()

Метод `.applymap()` позволяет применять функции с именованными аргументами ко всему датафрейму (метод `.apply()` применяется только к строкам или столбцам). Рассмотрим несложный пример.

```
1 # создадим датафрейм из чисел
2 nums_matrix = [[13, 7, 1],
3                 [4, 2, 25],
4                 [45, 3, 8]]
5
6 nums = pd.DataFrame(nums_matrix)
7 nums
```

	0	1	2
0	13	7	1
1	4	2	25
2	45	3	8

```
1 # объявим функцию, которая на входе принимает число x и
2 # прибавляет к нему другое число, указанное в параметре number
3 def add_number(x, number):
4     return x + number
5
6 # передадим методу .applymap() функцию add_number и
7 # прибавим единицу к каждому элементу датафрейма
8 nums.applymap(add_number, number = 1)
```

	0	1	2
0	14	8	2
1	5	3	26
2	46	4	9

Метод .pipe()

Метод `.pipe()`, как следует из его названия, позволяет создать pipeline и последовательно применить несколько функций к датафрейму. Вновь создадим исходный датафрейм с параметрами нескольких людей.

```
1 people = pd.DataFrame({'name'      : ['Алексей', 'Иван', 'Анна', 'Ольга', 'Николай'],
2                          'gender'   : [1, 1, 0, 2, 1],
3                          'age'      : [35, 20, 13, 28, 16],
4                          'height'   : [180.46, 182.26, 165.12, 168.04, 178.68],
5                          'weight'   : [73.61, 75.34, 50.22, 52.14, 69.72]
6                          })
7
8 people
```

	name	gender	age	height	weight
0	Алексей	1	35	180.46	73.61
1	Иван	1	20	182.26	75.34
2	Анна	0	13	165.12	50.22
3	Ольга	2	28	168.04	52.14
4	Николай	1	16	178.68	69.72

Объявим несколько функций, которые мы могли бы применить к датафрейму.

```
1 # в первую очередь скопируем датафрейм
2 def copy_df(df):
3     return df.copy()
4
5 # заменим значения столбца на новые с помощью метода .map()
6 def map_column(df, column, label1, label2):
7     labels_map = {0: label1, 1 : label2}
8     df[column] = df[column].map(labels_map)
9     return df
10
11 # кроме этого, создадим функцию для превращения количественной переменной
12 # в бинарную категориальную
13 def to_categorical(df, newcol, condcol, thres, cat1, cat2):
14     df[newcol] = np.where(df[condcol] >= thres, cat1, cat2)
15     return df
```

Последовательно применим эти функции с помощью нескольких методов `.pipe()`.

```
1 people_processed = (people.
2     pipe(copy_df). # copy_df() применится ко всему датафрейму
3     pipe(map_column, 'gender', 'female', 'male'). # map_column() к столбцу
4     pipe(to_categorical, 'age_group', 'age', 18, 'adult', 'minor')) # t
```

Посмотрим на результат и кроме того убедимся, что исходный датафрейм не изменился.

```
1 people_processed
```

	name	gender	age	height	weight	age_group
0	Алексей	male	35	180.46	73.61	adult
1	Иван	male	20	182.26	75.34	adult
2	Анна	female	13	165.12	50.22	minor
3	Ольга	NaN	28	168.04	52.14	adult
4	Николай	male	16	178.68	69.72	minor

```
1 people
```

	name	gender	age	height	weight
0	Алексей	1	35	180.46	73.61
1	Иван	1	20	182.26	75.34
2	Анна	0	13	165.12	50.22
3	Ольга	2	28	168.04	52.14
4	Николай	1	16	178.68	69.72

Перейдем к следующему разделу, который посвящен соединению датафреймов.

Соединение датафреймов

Рассмотрим, как мы можем соединить два датафрейма с помощью функций/методов `pd.concat()`, `pd.merge()` и `.join()`. Начнем с функции `pd.concat()`.

pd.concat()

В качестве примера возьмем информацию о стоимости канцелярских товаров в двух магазинах.

```
1 s1 = pd.DataFrame({
2     'item': ['карандаш', 'ручка', 'папка', 'степлер'],
3     'price': [220, 340, 200, 500]
4 })
5
6 s2 = pd.DataFrame({
7     'item': ['клей', 'корректор', 'скрепка', 'бумага'],
8     'price': [200, 240, 100, 300]
9 })
```

Выведем результат.

```
1 s1
```

	item	price
0	карандаш	220
1	ручка	340
2	папка	200
3	степлер	500

```
1 s2
```

	item	price
0	клей	200
1	корректор	240
2	скрепка	100
3	бумага	300

Соединение «один на другой»

В первую очередь мы можем совместить два датафрейма, поставив их «один на другой».

```
1 # передадим в функцию pd.concat() список из соединяемых датафреймов,
2 # укажем параметр axis = 0 (значение по умолчанию)
3 pd.concat([s1, s2], axis = 0)
```

	item	price
0	карандаш	220
1	ручка	340
2	папка	200
3	степлер	500
0	клей	200
1	корректор	240
2	скрепка	100
3	бумага	300

Как вы видите, индекс не обновился.

```
1 # обновим индекс через параметр ignore_index = True
2 pd.concat([s1, s2], axis = 0, ignore_index = True)
```

	item	price
0	карандаш	220
1	ручка	340
2	папка	200
3	степлер	500
4	клей	200
5	корректор	240
6	скрепка	100
7	бумага	300

Именно таким способом предлагается добавлять новые строки вместо удаленного из библиотеки метода `.append()`. При этом здесь есть нюанс, в функцию `pd.concat()` нельзя передать словарь, только объекты Series и DataFrame.

При соединении датафреймов мы можем создать многоуровневый (иерархический) индекс. Например, создадим отдельную группу для товаров первого магазина (s1) и товаров второго (s2).

```
1 # передадим в параметр keys названия групп индекса,
2 # параметр names получим названия уровней индекса
3 by_shop = pd.concat([s1, s2], axis = 0, keys = ['s1', 's2'], names = ['s', 'id'])
4 by_shop
```

		item	price
s	id		
s1	0	карандаш	220
	1	ручка	340
	2	папка	200
	3	степлер	500
s2	0	клей	200
	1	корректор	240
	2	скрепка	100
	3	бумага	300

Посмотрим на созданный индекс.

```
1 by_shop.index
```

```
1 MultiIndex([('s1', 0),
2             ('s1', 1),
3             ('s1', 2),
4             ('s1', 3),
5             ('s2', 0),
6             ('s2', 1),
7             ('s2', 2),
```

```

8 |         ('s2', 3)],
9 |         names=['s', 'id'])

```

Выведем первую запись в первой группе.

```

1 | by_shop.loc(['s1', 0])

```

```

1 | item      карандаш
2 | price      220
3 | Name: (s1, 0), dtype: object

```

Соединение «рядом друг с другом»

Датафреймы можно расположить «рядом друг с другом».

```

1 | # для этого сразу используем параметр axis = 1
2 | # одновременно сразу создадим группы для многоуровневого индекса столбцов
3 | pd.concat([s1, s2], axis = 1, keys = ['s1', 's2'])

```

	s1		s2	
	item	price	item	price
0	карандаш	220	клей	200
1	ручка	340	корректор	240
2	папка	200	скрепка	100
3	степлер	500	бумага	300

Выберем вторую группу (второй магазин) с помощью метода `.iloc[]`.

```

1 | pd.concat([s1, s2], axis = 1, keys = ['s1', 's2']).loc[:, 's2']

```

	item	price
0	клей	200
1	корректор	240
2	скрепка	100
3	бумага	300

Полученный через соединение

результат и в целом любой датафрейм можно транспонировать.

```

1 | # для транспонирования датафрейма используется метод .T или .transpose()
2 | pd.concat([s1, s2], axis = 1, keys = ['s1', 's2']).T

```

		0	1	2	3
s1	item	карандаш	ручка	папка	степлер
	price	220	340	200	500
s2	item	клей	корректор	скрепка	бумага
	price	200	240	100	300

Итак, `pd.concat()` выполняет простое «склеивание» датафреймов по вертикали или по горизонтали. Теперь посмотрим, что делать, если датафреймы нужно соединить по