

# Лекция 2: библиотека Pandas

**Автор: Сергей Вячеславович Макрушин** e-mail: SVMakrushin@fa.ru

Финансовый университет, 2020 г.

При подготовке лекции использованы материалы:

- Уэс Маккинли Python и анализ данных / Пер. с англ. Слипкин А.А. - М.: ДМК Пресс, 2015

V 0.4 10.09.2020

## Разделы:

- [Серии \(Series\) - одномерные массивы в Pandas](#)
- [Датафрэйм \(DataFrame\) - двумерные массивы в Pandas](#)
  - [Введение](#)
  - [Индексация](#)
- [Обработка данных в библиотеке Pandas](#)
  - [Универсальные функции и выравнивание](#)
  - [Работа с пустыми значениями](#)
  - [Агрегирование и группировка](#)
- [Обработка нескольких наборов данных](#)
  - [Объединение наборов данных](#)
  - [GroupBy: разбиение, применение, объединение](#)

-

- [к оглавлению](#)

```
In [2]: # загружаем стиль для оформления презентации
from IPython.display import HTML
from urllib.request import urlopen
html = urlopen("file:./lec_v1.css")
HTML(html.read().decode('utf-8'))
```

Out[2]:

**Pandas** - надстройка над библиотекой NumPy, обеспечивающая удобную инфраструктуру для обработки панельных данных (Pandas - от panel data sets).

Основным классом Pandas является **DataFrame**, объекты DataFrame - многомерные массивы с метками для строк и столбцов. DataFrame позволяет хранить:

- разнородные данные в различных столбцах
- корректно работать с пропущенными данными.

Кроме операций, поддерживаемых NumPy, библиотека Pandas реализует множество операций для работы с данными, характерных для работы с электронными таблицами и базами данных.

## Серии (Series) - одномерные массивы в Pandas

- [к оглавлению](#)

```
In [3]: import numpy as np
import pandas as pd
```

Фундаментальные структуры данных Pandas - классы **Series**, **DataFrame** и **Index**.

Объект **Series** - одномерный массив индексированных данных.

```
In [4]: # создание Series на основе списка Python:
srl = pd.Series([5, 6, 2, 9, 12])
srl
```

```
Out[4]: 0      5
1      6
2      2
3      9
4     12
dtype: int64
```

```
In [5]: srl.values # атрибут values - это массив NumPy со значениями
```

```
Out[5]: array([ 5,  6,  2,  9, 12], dtype=int64)
```

```
In [5]: srl.index # index - массивоподобный объект типа pd.Index
```

```
Out[5]: RangeIndex(start=0, stop=5, step=1)
```

```
In [6]: # Обращение к элементу серии по индексу:
srl[2]
```

```
Out[6]: 2
```

```
In [7]: # Срез серии по индексу:
srl[:3]
```

```
Out[7]: 0      5
1      6
2      2
dtype: int64
```

Основное различие между одномерным массивом библиотеки NumPy и Series - наличие у Series индекса, определяющего доступ к данным массива.

Индекс массива NumPy:

- всегда целочисленный
- представлен последовательно идущими целыми числами начиная с 0
- описывается неявно (т.е. не подразумевается явное определение индекса т.к. не допускаются альтернативные варианты индексации)

Индекс объекта Series:

- может состоять из значений типа, выбранного разработчиком (например, строк)

- индекс может описываться явно (вариант по умолчанию совпадает со способом индексации в NumPy) и связывается со значениями

```
In [8]: # Создание серии с явным определением индекса:
sr2 = pd.Series([5, 6, 2, 9, 12], index=['Cochise County', 'Pima County', 'Santa Cruz Cour
        'Maricopa County', 'Yuma County'])

sr2
```

```
Out[8]: Cochise County      5
Pima County      6
Santa Cruz County  2
Maricopa County  9
Yuma County     12
dtype: int64
```

```
In [9]: # Обращение к элементу серии по нецелочисленному индексу:
sr2['Pima County']
```

```
Out[9]: 6
```

```
In [10]: sr2['Pima County']
```

```
Out[10]: Pima County      6
Santa Cruz County  2
Maricopa County  9
Yuma County     12
dtype: int64
```

Объект Series можно рассматривать как специализированный вариант словаря.

- Словарь - структура, задающая соответствие произвольных ключей набору произвольных значений
- Объект Series:
  - структура, задающая соответствие **типизированных ключей** набору **типизированных значений**
  - кроме того, для ключей (значений индекса) задана **последовательность их следования**.

```
In [11]: # объект Series можно создавать непосредственно из словаря Python:
# (т.к. словарь не определяет порядок обхода, то такая форма задания может привести
# к созданию серии с иной последовательностью индекс-значение)
sr3 = pd.Series({'California': 38332521,
                 'Texas': 26448193,
                 'New York': 19651127,
                 'Florida': 19552860,
                 'Illinois': 12882135})

sr3
```

```
Out[11]: California      38332521
Texas      26448193
New York   19651127
Florida    19552860
Illinois   12882135
dtype: int64
```

```
In [12]: # изменение индекса:
sr3.index = ["Cochice", "Pima", "Santa Cruz", "Maricopa", "Yuma"]
sr3
```

```
Out[12]: Cochice      38332521
Pima      26448193
```

```
Santa Cruz      19651127
Maricopa        19552860
Yuma            12882135
dtype: int64
```

# Датафрэйм (DataFrame) - двумерные массивы в Pandas

- [к оглавлению](#)

## Введение

- [к оглавлению](#)

**DataFrame** - аналог двухмерного массива с гибкими индексами строк и гибкими именами столбцов.

Аналогично тому, что двумерный массив можно рассматривать как упорядоченную последовательность выровненных столбцов, объект DataFrame можно рассматривать как упорядоченную последовательность выровненных объектов Series. Под «выравниванием» понимается то, что они используют один и тот же индекс.

```
In [13]: # создание DataFrame на основе двух Series:
s_population = pd.Series({'California': 38332521,
                          'Texas': 26448193,
                          'New York': 19651127,
                          'Florida': 19552860,
                          'Illinois': 12882135})
s_area = pd.Series({'California': 423967, 'Texas': 695662, 'New York': 141297,
                   'Florida': 170312, 'Illinois': 149995})
states = pd.DataFrame({'population': s_population,
                       'area': s_area})
states # jupyter умеет красиво выводить таблицы Pandas DataFrame
```

```
Out[13]:
```

	population	area
<b>California</b>	38332521	423967
<b>Texas</b>	26448193	695662
<b>New York</b>	19651127	141297
<b>Florida</b>	19552860	170312
<b>Illinois</b>	12882135	149995

```
In [14]: # для всех столбцов DataFrame имеется единый индекс:
states.index
```

```
Out[14]: Index(['California', 'Texas', 'New York', 'Florida', 'Illinois'], dtype='object')
```

```
In [15]: # у объекта DataFrame есть атрибут columns, содержащий метки столбцов, - объект типа Index
states.columns
```

```
Out[15]: Index(['population', 'area'], dtype='object')
```

```
In [16]: # DataFrame можно рассматривать как специализированный словарь столбцов.
# DataFrame задает соответствие имени столбца объекту Series:
states['area']
```

```
Out[16]: California    423967
Texas          695662
New York       141297
Florida        170312
Illinois       149995
Name: area, dtype: int64
```

**NB!** Важно понимать, что в NumPy элементы по оси 0 принято рассматривать как **строки** (т.е. считается, что `np1[1]` - вернет строку с индексом 1), тогда как в Pandas аналогичная конструкция ( `pd1[1]` ) возвращает **столбец** типа Series.

```
In [17]: np1 = np.array([[1, 2, 3], [4, 5, 6]])
np1, np1.shape
```

```
Out[17]: (array([[1, 2, 3],
                [4, 5, 6]]),
         (2, 3))
```

```
In [18]: np1[1] # строка с индексом 1
```

```
Out[18]: array([4, 5, 6])
```

```
In [19]: # первое измерение (axis=0) рассматривается как размерность серий (столбцов),
# а вторая - как их количес
pd1 = pd.DataFrame(data=np1)
pd1
```

```
Out[19]:    0  1  2
0    1  2  3
1    4  5  6
```

```
In [20]: pd1[1] # обращение к столбцу с именем (индексом) 1
```

```
Out[20]: 0    2
1    5
Name: 1, dtype: int32
```

Т.е. индексация DataFrame (т.е. операция вида: `pd1[...]` ) ориентирована на манипулирование столбцами. *DataFrame* можно рассматривать как **серию серий** :

```
In [21]: type(pd1[1])
```

```
Out[21]: pandas.core.series.Series
```

```
In [22]: # из этого понятно, почему:
pd1[1][0]
```

```
Out[22]: 2
```

```
In [23]: # тогда как:
np1[1][0]
```

```
Out[23]: 4
```

```
In [24]: # создание DataFrame на базе массива NumPy с заданием индекса и имен столбцов
pd2 = pd.DataFrame(data=np1, index=['la', 'lb'], columns=['c11', 'c12', 'c13'] )
pd2
```

```
Out[24]:
```

	c11	c12	c13
la	1	2	3
lb	4	5	6

```
In [25]: # использование заданных индексов:
pd2['c12']
```

```
Out[25]:
```

la	2
lb	5

Name: c12, dtype: int32

```
In [26]: pd2['c12']['la']
```

```
Out[26]:
```

2

```
In [27]: # создание DataFrame из списка словарей (ключи - имена столбцов):
pd3 = pd.DataFrame([{'a': 1, 'b': 2, 'c': 'Alpha'}, {'a': 0, 'b': 3, 'c': 'Beta'}])
pd3
```

```
Out[27]:
```

	a	b	c
0	1	2	Alpha
1	0	3	Beta

```
In [28]: # явное задание индекса:
pd3 = pd.DataFrame([{'a': 1, 'b': 2, 'c': 'Alpha'}, {'a': 0, 'b': 3, 'c': 'Beta'}], index=['first', 'second'])
pd3
```

```
Out[28]:
```

	a	b	c
first	1	2	Alpha
second	0	3	Beta

```
In [29]: # в Pandas допускаются пропуски данных
# (и явная индексация упрощает задание данных с пропусками):
pd3 = pd.DataFrame([{'a': 1, 'c': 'Alpha'}, {'a': 0, 'b': 3, 'c': 'Beta'}], index=['first', 'second'])
pd3
```

```
Out[29]:
```

	a	c	b
first	1	Alpha	NaN
second	0	Beta	3.0

```
In [30]: # создание DataFrame из словаря списков (ключи - имена столбцов):
data = {'county': ['Cochice', 'Pima', 'Santa Cruz', 'Maricopa', 'Yuma'],
        'year': [2012, 2012, 2013, 2014, 2014],
        'reports': [4, 24, 31, 2, 3]}
```

```
pd4 = pd.DataFrame(data)
pd4
```

```
Out[30]:
```

	county	year	reports
0	Cochice	2012	4
1	Pima	2012	24
2	Santa Cruz	2013	31
3	Maricopa	2014	2
4	Yuma	2014	3

```
In [31]: # явное определение порядка и состава столбцов и индекса:
pd4 = pd.DataFrame(data, columns=['reports', 'county'], index=[chr(ord('a') + i) for i in
pd4
```

```
Out[31]:
```

	reports	county
a	4	Cochice
b	24	Pima
c	31	Santa Cruz
d	2	Maricopa
e	3	Yuma

## Индексация

- [к оглавлению](#)

### Индексация для серий

```
In [33]: sr4 = pd.Series([0.25, 0.5, 0.75, 1.0],
                        index=['a', 'b', 'c', 'd'])
sr4
```

```
Out[33]:
```

a	0.25
b	0.50
c	0.75
d	1.00

dtype: float64

Серии поддерживают интерфейс, близкий к словарям Python

```
In [34]: # извлечение элемента серии по аналогии с использованием словаря:
sr4['b']
```

```
Out[34]:
```

0.5

```
In [35]: # аналогично словарям поддерживается проверка вхождения элемента в индекс серии:
'a' in sr4
```

```
Out[35]:
```

True

```
In [36]:
```

```
sr4.keys()
```

```
Out[36]: Index(['a', 'b', 'c', 'd'], dtype='object')
```

```
In [37]: # в отличие от словарей keys() нужно указывать явно:
for i in sr4.keys():
    print(f'{i} -> {sr4[i]}')
```

```
a -> 0.25
b -> 0.5
c -> 0.75
d -> 1.0
```

```
In [38]: # итерация по значениям, а не по ключам!
for i in sr4:
    print(f'{i}')
```

```
0.25
0.5
0.75
1.0
```

```
In [39]: list(sr4.items())
```

```
Out[39]: [('a', 0.25), ('b', 0.5), ('c', 0.75), ('d', 1.0)]
```

```
In [40]: for i, v in sr4.items():
          print(f'{i} -> {v}')
```

```
a -> 0.25
b -> 0.5
c -> 0.75
d -> 1.0
```

```
In [41]: # модификация (добавление) элемента серии:
sr4['e'] = 1.25
sr4
```

```
Out[41]: a    0.25
         b    0.50
         c    0.75
         d    1.00
         e    1.25
         dtype: float64
```

```
In [42]: sr4['e'] = 1.75
sr4
```

```
Out[42]: a    0.25
         b    0.50
         c    0.75
         d    1.00
         e    1.75
         dtype: float64
```

Серии поддерживают механизмы индексации, аналогичные массивам NumPy: срезы, маскирование и прихотливое индексирование.

```
In [56]: # срез с использованием явных индексов (в срезах с явными использованием индексов правая л
```



```
sr4['a':'c']
```

```
Out[56]: a    0.25  
b    0.50  
c    0.75  
dtype: float64
```

```
In [44]: # прихотливое индексирование с использованием явных индексов:  
sr4[['b', 'a', 'c']]
```

```
Out[44]: b    0.50  
a    0.25  
c    0.75  
dtype: float64
```

```
In [45]: # срез с использованием НЕявных (целочисленных) индексов:  
sr4[0:2]
```

```
Out[45]: a    0.25  
b    0.50  
dtype: float64
```

```
In [46]: # прихотливое индексирование с использованием НЕявных индексов:  
sr4[[1, 0, 2]]
```

```
Out[46]: b    0.50  
a    0.25  
c    0.75  
dtype: float64
```

**NB!** В случае использования **НЕявного целочисленного индекса** использование срезов может выглядеть неоднозначно и **приводить к ошибкам**.

```
In [47]: sr5 = pd.Series(['a', 'b', 'c'], index=[1, 3, 5])
```

```
In [48]: # при обычном индексировании используется явный индекс  
sr5[1]
```

```
Out[48]: 'a'
```

```
In [49]: # при использовании среза используется НЕявный индекс:  
sr5[1:3] # этот результат может противоречить ожидаемому
```

```
Out[49]: 3    b  
5    c  
dtype: object
```

Из-за этой потенциальной путаницы в случае целочисленных индексов в библиотеке Pandas предусмотрены специальные атрибуты-индексаторы, позволяющие явным образом применять определенные схемы индексации:

- атрибут **loc** позволяет выполнить индексацию и срезы с использованием явного индекса
- атрибут **iloc** дает возможность выполнить индексацию и срезы, применяя неявный индекс в стиле языка Python

```
In [51]: sr5
```

```
Out[51]: 1    a
          3    b
          5    c
          dtype: object
```

```
In [52]: sr5.loc[1] # явный индекс
```

```
Out[52]: 'a'
```

```
In [53]: sr5.iloc[1] # неявный индекс
```

```
Out[53]: 'b'
```

```
In [57]: sr5.loc[1:3] # в срезах с явными использованием индексов правая граница включается!
```

```
Out[57]: 1    a
          3    b
          dtype: object
```

```
In [58]: sr5.iloc[1:3]
```

```
Out[58]: 3    b
          5    c
          dtype: object
```

---

```
In [59]: # Применение маскирования для серий аналогично NumPy:
          sr4[(sr4 > 0.3) & (sr4 < 0.8)]
```

```
Out[59]: b    0.50
          c    0.75
          dtype: float64
```

Что происходит внутри:

```
In [60]: sr4 > 0.3
```

```
Out[60]: a    False
          b     True
          c     True
          d     True
          e     True
          dtype: bool
```

```
In [61]: (sr4 > 0.3) & (sr4 < 0.8)
```

```
Out[61]: a    False
          b     True
          c     True
          d    False
          e    False
          dtype: bool
```

## Индексация для DataFrame

```
In [62]: states
```

```
Out[62]:
```

	population	area
<b>California</b>	38332521	423967
<b>Texas</b>	26448193	695662
<b>New York</b>	19651127	141297
<b>Florida</b>	19552860	170312
<b>Illinois</b>	12882135	149995

DataFrame может рассматриваться как словарь (серия) серий:

In [63]: `states['area']`

Out[63]:

California	423967
Texas	695662
New York	141297
Florida	170312
Illinois	149995

Name: area, dtype: int64

In [64]: *# для имен столбцов, не конфликтующих с методами DataFrame и синтаксисом Python, допустим*  
`states.area`

Out[64]:

California	423967
Texas	695662
New York	141297
Florida	170312
Illinois	149995

Name: area, dtype: int64

In [65]: *# синтаксис словаря допустим и для присвоения (создания новой серии-столбца):*  
`states['density'] = states['population'] / states['area']`  
`states`

Out[65]:

	population	area	density
<b>California</b>	38332521	423967	90.413926
<b>Texas</b>	26448193	695662	38.018740
<b>New York</b>	19651127	141297	139.076746
<b>Florida</b>	19552860	170312	114.806121
<b>Illinois</b>	12882135	149995	85.883763

**NB!** Операции среза и маскирования **относятся к строкам (!)**, а не столбцам (это не очень логично, но удобно на практике):

In [66]: `states[:'New York']` *# при явном использовании индекса правая граница включается!*

Out[66]:

	population	area	density
<b>California</b>	38332521	423967	90.413926
<b>Texas</b>	26448193	695662	38.018740
<b>New York</b>	19651127	141297	139.076746

```
In [67]: states[:3] # при НЕявном использовании индекса граница не включается
```

Out[67]:

	population	area	density
California	38332521	423967	90.413926
Texas	26448193	695662	38.018740
New York	19651127	141297	139.076746

```
In [68]: # маскирование работает по строкам:
states[states.density > 100]
```

Out[68]:

	population	area	density
New York	19651127	141297	139.076746
Florida	19552860	170312	114.806121

DataFrame поддерживает двухмерный вариант loc, iloc

```
In [69]: states.loc[states.density > 100, ['population', 'density']]
```

Out[69]:

	population	density
New York	19651127	139.076746
Florida	19552860	114.806121

```
In [70]: states.iloc[0, 2] = 90
states
```

Out[70]:

	population	area	density
California	38332521	423967	90.000000
Texas	26448193	695662	38.018740
New York	19651127	141297	139.076746
Florida	19552860	170312	114.806121
Illinois	12882135	149995	85.883763

# Обработка данных в библиотеке Pandas

- [к оглавлению](#)

## Универсальные функции и выравнивание

- [к оглавлению](#)

Все универсальные функции библиотеки NumPy работают с объектами Series и DataFrame библиотеки Pandas.

```
In [71]: import numpy as np
```

```
In [72]: rs = np.random.RandomState(42)
sr6 = pd.Series(rs.randint(0, 10, 4))
sr6
```

```
Out[72]: 0    6
1    3
2    7
3    4
dtype: int32
```

Результатом применения универсальной функции NumPy к объектам Pandas будет новый объект с сохранением индексов

```
In [73]: sr7 = np.exp(sr6)
sr7
```

```
Out[73]: 0    403.428793
1    20.085537
2   1096.633158
3    54.598150
dtype: float64
```

```
In [74]: sr6 # исходная серия осталась неизменной
```

```
Out[74]: 0    6
1    3
2    7
3    4
dtype: int32
```

```
In [75]: pd5 = pd.DataFrame(rs.randint(0, 10, (3, 4)),
                             columns=['A', 'B', 'C', 'D'])
pd5
```

```
Out[75]:
```

	A	B	C	D
0	6	9	2	6
1	7	4	3	7
2	7	2	5	4

```
In [76]: np.sin(pd5 * np.pi / 4)
```

```
Out[76]:
```

	A	B	C	D
0	-1.000000	7.071068e-01	1.000000	-1.000000e+00
1	-0.707107	1.224647e-16	0.707107	-7.071068e-01
2	-0.707107	1.000000e+00	-0.707107	1.224647e-16

При бинарных операциях над двумя объектами Series или DataFrame библиотека Pandas будет выравнивать индексы в процессе выполнения операции. Получившийся в итоге массив содержит объединение индексов двух исходных массивов. Недостающие значения будут отмечены как NaN («нечисловое значение»), с помощью которого библиотека Pandas отмечает пропущенные данные.

```
In [77]: pd5
```

Out[77]:

	A	B	C	D
0	6	9	2	6
1	7	4	3	7
2	7	2	5	4

```
In [78]: pd6 = pd.DataFrame(rs.randint(0, 10, (4, 4)), index=list(range(1,5)),
                        columns=['B', 'C', 'D', 'F'])
pd6
```

Out[78]:

	B	C	D	F
1	1	7	5	1
2	4	0	9	5
3	8	0	9	2
4	6	3	8	2

```
In [79]: sr8 = pd5['A'] + pd6['B'] # выполняется выравнивание по индексам (участвуют две серии)
sr8
```

Out[79]:

0	NaN
1	8.0
2	11.0
3	NaN
4	NaN

dtype: float64

```
In [80]: pd7 = pd5 + pd6 # выполняется выравнивание по столбцам и по индексам
pd7
```

Out[80]:

	A	B	C	D	F
0	NaN	NaN	NaN	NaN	NaN
1	NaN	5.0	10.0	12.0	NaN
2	NaN	6.0	5.0	13.0	NaN
3	NaN	NaN	NaN	NaN	NaN
4	NaN	NaN	NaN	NaN	NaN

## Работа с пустыми значениями

- [к оглавлению](#)

В Pandas в качестве пустых значений рассматривается значение NaN ("Not a Number"), поддерживаемое форматом чисел с плавающей точкой ( np.nan в NumPy) и значением None для объектов Python.

```
In [81]: sr8
```

Out[81]:

0	NaN
1	8.0
2	11.0
3	NaN

```
4      NaN
dtype: float64
```

```
In [82]: # получение маски пустых значений
sr8.isna()
```

```
Out[82]: 0      True
1      False
2      False
3       True
4       True
dtype: bool
```

```
In [83]: pd7.isna()
```

```
Out[83]:
```

	A	B	C	D	F
0	True	True	True	True	True
1	True	False	False	False	True
2	True	False	False	False	True
3	True	True	True	True	True
4	True	True	True	True	True

```
In [84]: # очистка от пустых значений:
sr8.dropna()
```

```
Out[84]: 1      8.0
2     11.0
dtype: float64
```

```
In [85]: pd7.dropna() # default how='any'
```

```
Out[85]:
```

	A	B	C	D	F
--	---	---	---	---	---

```
In [86]: # default axis=0, удаляем строки, в которых все значения NaN:
pd7.dropna(how='all')
```

```
Out[86]:
```

	A	B	C	D	F
1	NaN	5.0	10.0	12.0	NaN
2	NaN	6.0	5.0	13.0	NaN

```
In [87]: # последовательное применение dropna:
# сначала для строк (т.к. default axis=0),
# потом для столбцов dropna(axis=1), помним: (default how='any'):
pd7.dropna(how='all').dropna(axis=1)
```

```
Out[87]:
```

	B	C	D
1	5.0	10.0	12.0
2	6.0	5.0	13.0

In [88]: `pd7.fillna(0.0) # заполнение NaN заданными значениями`

Out[88]:

	A	B	C	D	F
0	0.0	0.0	0.0	0.0	0.0
1	0.0	5.0	10.0	12.0	0.0
2	0.0	6.0	5.0	13.0	0.0
3	0.0	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0	0.0	0.0

## Агрегирование и группировка

- [к оглавлению](#)

In [90]: `pd75 = pd.DataFrame({'A': rs.rand(5), 'B': rs.rand(5)})`  
`pd75`

Out[90]:

	A	B
0	0.683264	0.182236
1	0.609997	0.755361
2	0.833195	0.425156
3	0.173365	0.207942
4	0.391061	0.567700

In [91]: `pd75.values.mean(axis=0)`

Out[91]: `array([0.53817607, 0.42767907])`

In [92]: `# default axis=0, т.е. агрегируем значения вдоль оси 0`  
`# (т.е. при агрегировании меняем индекс элементов вдоль этой оси):`  
`pd75.mean()`

Out[92]:

A	0.538176
B	0.427679

`dtype: float64`

In [93]: `pd75.mean(axis=1)`

Out[93]:

0	0.432750
1	0.682679
2	0.629175
3	0.190653
4	0.479380

`dtype: float64`

In [96]: `# агрегирование по всему DataFrame:`  
`pd75.values.mean()`

Out[96]: `0.48292757127103964`



```
In [97]: # атрибут values:
pd75.values, type(pd75.values)
```

```
Out[97]: (array([[0.68326352, 0.18223609],
        [0.60999666, 0.75536141],
        [0.83319491, 0.42515587],
        [0.17336465, 0.20794166],
        [0.39106061, 0.56770033]]),
numpy.ndarray)
```

```
In [98]: pd.DataFrame({'sum':pd75.sum(), 'prod':pd75.prod(),
        'mean':pd75.mean(), 'median':pd75.median(), 'std':pd75.std(), 'var':pd75.var(),
        'min':pd75.min(), 'max':pd75.max()})
```

```
Out[98]:
```

	sum	prod	mean	median	std	var	min	max
A	2.690880	0.023543	0.538176	0.609997	0.258832	0.066994	0.173365	0.833195
B	2.138395	0.006909	0.427679	0.425156	0.242649	0.058879	0.182236	0.755361

```
In [99]: pd75.describe()
```

```
Out[99]:
```

	A	B
count	5.000000	5.000000
mean	0.538176	0.427679
std	0.258832	0.242649
min	0.173365	0.182236
25%	0.391061	0.207942
50%	0.609997	0.425156
75%	0.683264	0.567700
max	0.833195	0.755361

```
In [100...: # квантиль:
pd75.quantile(0.5)
```

```
Out[100...: A    0.609997
B    0.425156
Name: 0.5, dtype: float64
```

```
In [101...: pd75.quantile(np.arange(0.0, 1.1, 0.1))
```

```
Out[101...:
```

	A	B
0.0	0.173365	0.182236
0.1	0.260443	0.192518
0.2	0.347521	0.202801
0.3	0.434848	0.251385
0.4	0.522422	0.338270
0.5	0.609997	0.425156

	A	B
0.6	0.639303	0.482174
0.7	0.668610	0.539191
0.8	0.713250	0.605233
0.9	0.773222	0.680297
1.0	0.833195	0.755361

## Обработка нескольких наборов данных

### Объединение наборов данных

- [к оглавлению](#)

```
In [102... pd8 = pd.DataFrame([[1, 2], [3, 4]], columns=list('AB'))
pd8
```

```
Out[102...   A  B
0  1  2
1  3  4
```

```
In [103... pd9 = pd.DataFrame([[5, 6], [7, 8]], columns=list('AB'))
pd9
```

```
Out[103...   A  B
0  5  6
1  7  8
```

```
In [104... # append создает новый объект DataFrame:
pd8.append(pd9) # при конкатенации может происходить дублирование индекса
```

```
Out[104...   A  B
0  1  2
1  3  4
0  5  6
1  7  8
```

```
In [105... # автоматически создается новый индекс:
pd8.append(pd9, ignore_index=True)
```

```
Out[105...   A  B
0  1  2
1  3  4
```

**A B**

**2** 5 6

**3** 7 8

Функция `pd.merge()` реализует множество типов соединений: «один-к-одному», «многие-к-одному» и «многие-ко-многим». Все эти три типа соединений доступны через один и тот же вызов `pd.merge()`, тип выполняемого соединения зависит от формы входных данных.

In [106...

```
pd10 = pd.DataFrame({'employee': ['Bob', 'Jake', 'Lisa', 'Sue'],  
                     'group': ['Accounting', 'Engineering', 'Engineering', 'HR']})  
pd11 = pd.DataFrame({'employee': ['Lisa', 'Bob', 'Jake', 'Sue'],  
                     'hire_date': [2004, 2008, 2012, 2014]})
```

In [107...

pd10

Out[107...

	employee	group
0	Bob	Accounting
1	Jake	Engineering
2	Lisa	Engineering
3	Sue	HR

In [108...

pd11

Out[108...

	employee	hire_date
0	Lisa	2004
1	Bob	2008
2	Jake	2012
3	Sue	2014

In [131...

```
# Функция pd.merge() распознает, что в обоих объектах DataFrame имеется столбец  
# employee, и автоматически выполняет соединение один-к-одному, используя этот столбец в  
pd.merge(pd10, pd11)
```

Out[131...

	employee	group	hire_date
0	Bob	Accounting	2008
1	Jake	Engineering	2012
2	Lisa	Engineering	2004
3	Sue	HR	2014

In [109...

```
pd12 = pd.DataFrame({'group': ['Accounting', 'Engineering', 'HR'],  
                     'supervisor': ['Carly', 'Guido', 'Steve']})  
pd12
```

Out[109...

	<b>group</b>	<b>supervisor</b>
<b>0</b>	Accounting	Carly
<b>1</b>	Engineering	Guido
<b>2</b>	HR	Steve

```
In [110... # соединение многие-к-одному по столбцу group:
pd.merge(pd10, pd12)
```

	<b>employee</b>	<b>group</b>	<b>supervisor</b>
<b>0</b>	Bob	Accounting	Carly
<b>1</b>	Jake	Engineering	Guido
<b>2</b>	Lisa	Engineering	Guido
<b>3</b>	Sue	HR	Steve

```
In [111... pd13 = pd.DataFrame({'group': ['Accounting', 'Accounting', 'Engineering', 'Engineering', 'HR', 'HR'],
                        'skills': ['math', 'spreadsheets', 'coding', 'linux', 'spreadsheets', 'organization']})
pd13
```

	<b>group</b>	<b>skills</b>
<b>0</b>	Accounting	math
<b>1</b>	Accounting	spreadsheets
<b>2</b>	Engineering	coding
<b>3</b>	Engineering	linux
<b>4</b>	HR	spreadsheets
<b>5</b>	HR	organization

```
In [112... # соединение многие-ко-многим по столбцу group:
pd.merge(pd10, pd13)
```

	<b>employee</b>	<b>group</b>	<b>skills</b>
<b>0</b>	Bob	Accounting	math
<b>1</b>	Bob	Accounting	spreadsheets
<b>2</b>	Jake	Engineering	coding
<b>3</b>	Jake	Engineering	linux
<b>4</b>	Lisa	Engineering	coding
<b>5</b>	Lisa	Engineering	linux
<b>6</b>	Sue	HR	spreadsheets
<b>7</b>	Sue	HR	organization

Метод `pd.merge()` по умолчанию выполняет поиск в двух входных объектах соответствующих названий столбцов и использует найденное в качестве ключа. Однако, зачастую имена столбцов не совпадают, для

этого случая в методе `pd.merge()` имеются специальные параметры.

- `on` для явного указания имени (имен) столбцов;
- `left_on` и `right_on` для явного указания имен столбцов, в случае, если у первого и второго `DataFrame` они не совпадают;
- `left_index` и `right_index` для указания индекса в качестве ключа слияния.

In [113...

```
# пример:
pd14 = pd.DataFrame({'name': ['Bob', 'Jake', 'Lisa', 'Sue'],
                     'salary': [70000, 80000, 120000, 90000]})
pd15 = pd.merge(pd10, pd14, left_on='employee', right_on='name')
pd15
```

Out[113...

	employee	group	name	salary
0	Bob	Accounting	Bob	70000
1	Jake	Engineering	Jake	80000
2	Lisa	Engineering	Lisa	120000
3	Sue	HR	Sue	90000

In [114...

```
# лишний столбец можно удалить:
pd15.drop('name', axis=1, inplace=True) # inplace=True - не создается новый DataFrame
pd15
```

Out[114...

	employee	group	salary
0	Bob	Accounting	70000
1	Jake	Engineering	80000
2	Lisa	Engineering	120000
3	Sue	HR	90000

## GroupBy: разбиение, применение, объединение

- [к оглавлению](#)

Операцию `GroupBy` удобно представить в виде последовательного применения операций: разбиение, применение и объединение (**split, apply, combine**):

- **split** (шаг разбиения): включает разделение на части и группировку объекта `DataFrame` на основе значений заданного ключа.
- **apply** (шаг применения): включает вычисление какой-либо функции, обычно агрегирующей, преобразование или фильтрацию в пределах отдельных групп.
- **combine** (шаг объединения): во время шага выполняется слияние результатов предыдущих операций в выходной массив.

Для `DataFrame` операцию "разбить, применить, объединить" можно реализовать с помощью метода `groupby()`, передав в него имя желаемого ключевого столбца. Функция `groupby()` возвращает не набор объектов `DataFrame`, а объект `DataFrameGroupBy`, который можно рассматривать как специальное представление объекта `DataFrame`, готовое к группировке, но не выполняющее никаких фактических вычислений до этапа применения агрегирования (используется принцип отложенных вычислений).

Для получения результата нужно вызвать один из агрегирующих методов объекта DataFrameGroupBy, что приведет к выполнению соответствующих шагов применения/объединения.

```
In [117... pdl6 = pd.DataFrame({'key': ['A', 'B', 'C', 'A', 'B', 'C'],  
                      'data': range(1, 7)}, columns=['key', 'data'])
```

```
In [118... pdl6
```

```
Out[118...   key  data  
0    A     1  
1    B     2  
2    C     3  
3    A     4  
4    B     5  
5    C     6
```

```
In [119... pdl6.groupby('key')
```

```
Out[119... <pandas.core.groupby.generic.DataFrameGroupBy object at 0x000002DEAD6A49C8>
```

```
In [120... pdl6.groupby('key').sum()
```

```
Out[120...   data  
key  
A      5  
B      7  
C      9
```

```
In [121... # загружаем набор данных об открытии экзопланет:  
import seaborn as sns  
planets = sns.load_dataset('planets')  
planets.shape
```

```
Out[121... (1035, 6)
```

```
In [122... # заголовок таблицы  
planets.head()
```

```
Out[122...   method  number  orbital_period  mass  distance  year  
0  Radial Velocity      1         269.300   7.10      77.40  2006  
1  Radial Velocity      1         874.774   2.21      56.95  2008  
2  Radial Velocity      1         763.000   2.60      19.84  2011  
3  Radial Velocity      1         326.030  19.40     110.62  2007
```

	method	number	orbital_period	mass	distance	year
4	Radial Velocity	1	516.220	10.50	119.47	2009

```
In [123...  
# подсчитываем количество не NaN значений в каждой группе:  
planets.groupby('year').count()
```

	method	number	orbital_period	mass	distance
year					
1989		1	1	1	1
1992		2	2	0	0
1994		1	1	0	0
1995		1	1	1	1
1996		6	6	4	6
1997		1	1	1	1
1998		5	5	5	5
1999		15	15	14	15
2000		16	16	14	16
2001		12	12	11	12
2002		32	32	31	31
2003		25	25	22	24
2004		26	26	22	15
2005		39	39	38	34
2006		31	31	28	20
2007		53	53	52	32
2008		74	74	69	43
2009		98	98	96	74
2010		102	102	96	41
2011		185	185	184	91
2012		140	140	132	24
2013		118	118	107	30
2014		52	52	51	5

```
In [124...  
# группировка экзопланет по методу их идентификации:  
planets.groupby('method').count()
```

	number	orbital_period	mass	distance	year
method					
Astrometry	2	2	0	2	2
Eclipse Timing Variations	9	9	2	4	9

	number	orbital_period	mass	distance	year
<b>method</b>					
<b>Imaging</b>	38	12	0	32	38
<b>Microlensing</b>	23	7	0	10	23
<b>Orbital Brightness Modulation</b>	3	3	0	2	3
<b>Pulsar Timing</b>	5	5	0	1	5
<b>Pulsation Timing Variations</b>	1	1	0	0	1
<b>Radial Velocity</b>	553	553	510	530	553
<b>Transit</b>	397	397	1	224	397
<b>Transit Timing Variations</b>	4	3	0	3	4

```
In [125... # сколько орбитальных периодов было обнаружено каждым из методов:
planets.groupby('method')['orbital_period'].count()
```

```
Out[125... method
Astrometry          2
Eclipse Timing Variations  9
Imaging             12
Microlensing         7
Orbital Brightness Modulation  3
Pulsar Timing        5
Pulsation Timing Variations  1
Radial Velocity     553
Transit             397
Transit Timing Variations  3
Name: orbital_period, dtype: int64
```

```
In [126... # медианное значение орбитальных периодов (в днях), выявленных каждым из методов:
planets.groupby('method')['orbital_period'].median()
```

```
Out[126... method
Astrometry          631.180000
Eclipse Timing Variations  4343.500000
Imaging             27500.000000
Microlensing         3300.000000
Orbital Brightness Modulation  0.342887
Pulsar Timing        66.541900
Pulsation Timing Variations  1170.000000
Radial Velocity     360.200000
Transit             5.714932
Transit Timing Variations  57.011000
Name: orbital_period, dtype: float64
```

```
In [222... # по группам, выделенным с помощью groupby, можно итерироваться:
for (method, group) in planets.groupby('method'): # тип group - DataFrame
    print(f'{method} shape={group.shape}')
```

```
Astrometry shape=(2, 6)
Eclipse Timing Variations shape=(9, 6)
Imaging shape=(38, 6)
Microlensing shape=(23, 6)
Orbital Brightness Modulation shape=(3, 6)
Pulsar Timing shape=(5, 6)
Pulsation Timing Variations shape=(1, 6)
Radial Velocity shape=(553, 6)
```



```
Transit shape=(397, 6)
Transit Timing Variations shape=(4, 6)
```

На этапе применения у объектов GroupBy кроме обычных агрегирующих методов, таких как `sum()`, `median()` и т. п., имеются методы `aggregate()`, `filter()`, `transform()` и `apply()`, эффективно выполняющие множество полезных операций до объединения сгруппированных данных.

Метод `aggregate()` может принимать на входе строку, функцию или список и вычислять все сводные показатели сразу.

In [229...

```
planets.groupby('method')['orbital_period'].aggregate(['min', np.median, max])
```

Out[229...

	min	median	max
method			
Astrometry	246.360000	631.180000	1016.000000
Eclipse Timing Variations	1916.250000	4343.500000	10220.000000
Imaging	4639.150000	27500.000000	730000.000000
Microlensing	1825.000000	3300.000000	5100.000000
Orbital Brightness Modulation	0.240104	0.342887	1.544929
Pulsar Timing	0.090706	66.541900	36525.000000
Pulsation Timing Variations	1170.000000	1170.000000	1170.000000
Radial Velocity	0.736540	360.200000	17337.500000
Transit	0.355000	5.714932	331.600590
Transit Timing Variations	22.339500	57.011000	160.000000

Операция фильтрации `filter` дает возможность опускать данные в зависимости от свойств группы. Например, нам может понадобиться оставить в результате все группы

In [263...

```
def filter_func(x):
    return x['orbital_period'].max()/x['orbital_period'].min() > 1000
```

In [264...

```
gr1 = planets.groupby('method').filter(filter_func)
gr1.shape
```

Out[264...

```
(558, 6)
```

В то время как агрегирующая функция должна возвращать сокращенную версию данных, преобразование `transform` может вернуть версию полного набора данных, преобразованную ради дальнейшей их перекомпоновки. При подобном преобразовании форма выходных данных совпадает с формой входных. Распространенный пример - центрирование данных путем вычитания среднего значения по группам.

In [267...

```
planets['cntr_orbital_period'] = planets.groupby('method')['orbital_period'].transform(lambda x: x - x.mean())
```

Out[267...

	method	number	orbital_period	mass	distance	year	cntr_orbital_period
0	Radial Velocity	1	269.300000	7.100	77.40	2006	-554.054680

	method	number	orbital_period	mass	distance	year	cntr_orbital_period
1	Radial Velocity	1	874.774000	2.210	56.95	2008	51.419320
2	Radial Velocity	1	763.000000	2.600	19.84	2011	-60.354680
3	Radial Velocity	1	326.030000	19.400	110.62	2007	-497.324680
4	Radial Velocity	1	516.220000	10.500	119.47	2009	-307.134680
5	Radial Velocity	1	185.840000	4.800	76.39	2008	-637.514680
6	Radial Velocity	1	1773.400000	4.640	18.15	2002	950.045320
7	Radial Velocity	1	798.500000	NaN	21.41	1996	-24.854680
8	Radial Velocity	1	993.300000	10.300	73.10	2008	169.945320
9	Radial Velocity	2	452.800000	1.990	74.79	2010	-370.554680
10	Radial Velocity	2	883.000000	0.860	74.79	2010	59.645320
11	Radial Velocity	1	335.100000	9.880	39.43	2009	-488.254680
12	Radial Velocity	1	479.100000	3.880	97.28	2008	-344.254680
13	Radial Velocity	3	1078.000000	2.530	14.08	1996	254.645320
14	Radial Velocity	3	2391.000000	0.540	14.08	2001	1567.645320
15	Radial Velocity	3	14002.000000	1.640	14.08	2009	13178.645320
16	Radial Velocity	1	4.230785	0.472	15.36	1995	-819.123895
17	Radial Velocity	5	14.651000	0.800	12.53	1996	-808.703680
18	Radial Velocity	5	44.380000	0.165	12.53	2004	-778.974680
19	Radial Velocity	5	4909.000000	3.530	12.53	2002	4085.645320
20	Radial Velocity	5	0.736540	NaN	12.53	2011	-822.618140
21	Radial Velocity	5	261.200000	0.172	12.53	2007	-562.154680
22	Radial Velocity	3	4.215000	0.016	8.52	2009	-819.139680
23	Radial Velocity	3	38.021000	0.057	8.52	2009	-785.333680
24	Radial Velocity	3	123.010000	0.072	8.52	2009	-700.344680
25	Radial Velocity	1	116.688400	NaN	18.11	1996	-706.666280
26	Radial Velocity	1	691.900000	NaN	81.50	2012	-131.454680
27	Radial Velocity	1	952.700000	5.300	97.18	2008	129.345320
28	Radial Velocity	1	181.400000	3.200	45.52	2013	-641.954680
29	Imaging	1	NaN	NaN	45.52	2005	NaN
...	...	...	...	...	...	...	...
1005	Transit	1	3.693641	NaN	200.00	2012	-17.408432
1006	Transit	1	4.465633	NaN	330.00	2012	-16.636440
1007	Transit	1	4.617101	NaN	255.00	2012	-16.484972
1008	Transit	1	2.838971	NaN	455.00	2012	-18.263102
1009	Transit	1	5.017180	NaN	300.00	2012	-16.084893
1010	Transit	1	7.919585	NaN	125.00	2012	-13.182488

	method	number	orbital_period	mass	distance	year	cntr_orbital_period
1011	Transit	1	4.305001	NaN	400.00	2012	-16.797072
1012	Transit	1	3.855900	NaN	480.00	2012	-17.246173
1013	Transit	1	4.411953	NaN	160.00	2012	-16.690120
1014	Transit	1	4.378090	NaN	330.00	2012	-16.723983
1015	Transit	1	1.573292	NaN	350.00	2012	-19.528781
1016	Transit	1	2.311424	NaN	310.00	2013	-18.790648
1017	Transit	1	4.086052	NaN	380.00	2012	-17.016021
1018	Transit	1	4.614420	NaN	225.00	2012	-16.487653
1019	Transit	1	2.903675	NaN	345.00	2012	-18.198398
1020	Transit	1	2.216742	NaN	340.00	2012	-18.885331
1021	Transit	1	2.484193	NaN	260.00	2013	-18.617880
1022	Transit	1	1.360031	NaN	93.00	2012	-19.742042
1023	Transit	1	2.175176	NaN	550.00	2012	-18.926896
1024	Transit	1	3.662387	NaN	240.00	2012	-17.439686
1025	Transit	1	3.067850	NaN	60.00	2012	-18.034222
1026	Transit	1	0.925542	NaN	470.00	2014	-20.176531
1027	Imaging	1	NaN	NaN	19.20	2011	NaN
1028	Transit	1	3.352057	NaN	3200.00	2012	-17.750016
1029	Imaging	1	NaN	NaN	10.10	2012	NaN
1030	Transit	1	3.941507	NaN	172.00	2006	-17.160566
1031	Transit	1	2.615864	NaN	148.00	2007	-18.486209
1032	Transit	1	3.191524	NaN	174.00	2007	-17.910549
1033	Transit	1	4.125083	NaN	293.00	2008	-16.976990
1034	Transit	1	4.187757	NaN	260.00	2008	-16.914316

1035 rows × 7 columns

Метод `apply()` позволяет применять произвольную функцию к результатам группировки. В качестве параметра эта функция должна получать объект `DataFrame`, а возвращать или объект библиотеки `Pandas` (например, `DataFrame`, `Series`), или скалярное значение, в зависимости от возвращаемого значения будет вызвана соответствующая операция объединения.

In [146...

```
def norm_by_min_in_year(x):
    # x - объект DataFrame сгруппированных значений
    x['orbital_period_normalized'] = x['orbital_period']/x['orbital_period'].min()
    return x
```

In [147...

```
planets.groupby('year').apply(norm_by_min_in_year)
```

Out[147...

	method	number	orbital_period	mass	distance	year	orbital_period_normalized
--	--------	--------	----------------	------	----------	------	---------------------------

	method	number	orbital_period	mass	distance	year	orbital_period_normalized
0	Radial Velocity	1	269.300000	7.100	77.40	2006	149.944321
1	Radial Velocity	1	874.774000	2.210	56.95	2008	801.498594
2	Radial Velocity	1	763.000000	2.600	19.84	2011	8411.765050
3	Radial Velocity	1	326.030000	19.400	110.62	2007	249.604610
4	Radial Velocity	1	516.220000	10.500	119.47	2009	654.403935
5	Radial Velocity	1	185.840000	4.800	76.39	2008	170.273121
6	Radial Velocity	1	1773.400000	4.640	18.15	2002	1463.299236
7	Radial Velocity	1	798.500000	NaN	21.41	1996	240.983854
8	Radial Velocity	1	993.300000	10.300	73.10	2008	910.096269
9	Radial Velocity	2	452.800000	1.990	74.79	2010	373.325067
10	Radial Velocity	2	883.000000	0.860	74.79	2010	728.016859
11	Radial Velocity	1	335.100000	9.880	39.43	2009	424.800974
12	Radial Velocity	1	479.100000	3.880	97.28	2008	438.968209
13	Radial Velocity	3	1078.000000	2.530	14.08	1996	325.335748
14	Radial Velocity	3	2391.000000	0.540	14.08	2001	380.975143
15	Radial Velocity	3	14002.000000	1.640	14.08	2009	17750.114092
16	Radial Velocity	1	4.230785	0.472	15.36	1995	1.000000
17	Radial Velocity	5	14.651000	0.800	12.53	1996	4.421609
18	Radial Velocity	5	44.380000	0.165	12.53	2004	30.981339
19	Radial Velocity	5	4909.000000	3.530	12.53	2002	4050.601076
20	Radial Velocity	5	0.736540	NaN	12.53	2011	8.120054
21	Radial Velocity	5	261.200000	0.172	12.53	2007	199.971549
22	Radial Velocity	3	4.215000	0.016	8.52	2009	5.343289
23	Radial Velocity	3	38.021000	0.057	8.52	2009	48.198621
24	Radial Velocity	3	123.010000	0.072	8.52	2009	155.937833
25	Radial Velocity	1	116.688400	NaN	18.11	1996	35.216056
26	Radial Velocity	1	691.900000	NaN	81.50	2012	931.274900
27	Radial Velocity	1	952.700000	5.300	97.18	2008	872.897126
28	Radial Velocity	1	181.400000	3.200	45.52	2013	510.985915
29	Imaging	1	NaN	NaN	45.52	2005	NaN
...	...	...	...	...	...	...	...
1005	Transit	1	3.693641	NaN	200.00	2012	4.971521
1006	Transit	1	4.465633	NaN	330.00	2012	6.010597
1007	Transit	1	4.617101	NaN	255.00	2012	6.214468
1008	Transit	1	2.838971	NaN	455.00	2012	3.821163
1009	Transit	1	5.017180	NaN	300.00	2012	6.752961

	method	number	orbital_period	mass	distance	year	orbital_period_normalized
1010	Transit	1	7.919585	NaN	125.00	2012	10.659504
1011	Transit	1	4.305001	NaN	400.00	2012	5.794391
1012	Transit	1	3.855900	NaN	480.00	2012	5.189916
1013	Transit	1	4.411953	NaN	160.00	2012	5.938345
1014	Transit	1	4.378090	NaN	330.00	2012	5.892767
1015	Transit	1	1.573292	NaN	350.00	2012	2.117600
1016	Transit	1	2.311424	NaN	310.00	2013	6.511054
1017	Transit	1	4.086052	NaN	380.00	2012	5.499693
1018	Transit	1	4.614420	NaN	225.00	2012	6.210859
1019	Transit	1	2.903675	NaN	345.00	2012	3.908252
1020	Transit	1	2.216742	NaN	340.00	2012	2.983663
1021	Transit	1	2.484193	NaN	260.00	2013	6.997727
1022	Transit	1	1.360031	NaN	93.00	2012	1.830557
1023	Transit	1	2.175176	NaN	550.00	2012	2.927717
1024	Transit	1	3.662387	NaN	240.00	2012	4.929453
1025	Transit	1	3.067850	NaN	60.00	2012	4.129227
1026	Transit	1	0.925542	NaN	470.00	2014	1.382830
1027	Imaging	1	NaN	NaN	19.20	2011	NaN
1028	Transit	1	3.352057	NaN	3200.00	2012	4.511760
1029	Imaging	1	NaN	NaN	10.10	2012	NaN
1030	Transit	1	3.941507	NaN	172.00	2006	2.194603
1031	Transit	1	2.615864	NaN	148.00	2007	2.002674
1032	Transit	1	3.191524	NaN	174.00	2007	2.443392
1033	Transit	1	4.125083	NaN	293.00	2008	3.779545
1034	Transit	1	4.187757	NaN	260.00	2008	3.836970

1035 rows × 7 columns

## Технический раздел

next **Q:** qs line

next **A:** an line

next **Def:** df line

next **Ex:** ex line

next **+** pl line

next **-** mn line

next  $\pm$  plmn line

next  $\Rightarrow$  hn line

$\Rightarrow$  Home

- News red and green and blue and \_\_selected\_\_
- $A \Rightarrow b \Rightarrow c \triangleright$  Contact
- $\approx \sim \approx \pm$  About

- **Def:** Определение
- **Ex:** пример (кейс)
- **Q:** вопрос (проблема)
- **A:** ответ
- Алгоритм:
  - **S1:** Шаг 1
  - **S2:** Шаг 2
- Свойства:
  - **P1:** Свойство 1
  - **P2:** Свойство 2
- Утверждение
  - $\Rightarrow$  следствие
- Свойства:
  - + положительные
  - - отрицательные
  - $\pm$  смешанные