

Aggregation and Grouping

聚合与分组

An essential piece of analysis of large data is efficient summarization: computing aggregations like `sum()`, `mean()`, `median()`, `min()`, and `max()`, in which a single number gives insight into the nature of a potentially large dataset. In this section, we'll explore aggregations in Pandas, from simple operations akin to what we've seen on NumPy arrays, to more sophisticated operations based on the concept of a "groupby".

对于一个大数据集进行分析的使用有效的概括: 对数据集进行 `sum()`、`mean()`、`median()`、`min()` 和 `max()` 聚合运算, 这些运算的结果就可能可以给出大数据集的一些内在特征。在本节中, 我们会探讨Pandas中的聚合, 从我们已经在NumPy数组中进行的这些简单的操作, 直到基于分组 `groupby` 概念进行的更复杂的操作。

For convenience, we'll use the same `display` magic function that we've seen in previous sections:

```
In [1]: import numpy as np
import pandas as pd

class display(object):
    """Display HTML representation of multiple objects"""
    template = """<div style="float: left; padding: 10px;">
<p style="font-family:'Courier New', Courier, monospace">{0}</p>{1}
</div>"""
    def __init__(self, *args):
        self.args = args

    def repr_html(self):
        return '\n'.join(self.template.format(a, eval(a)._repr_html_())
                        for a in self.args)

    def __repr__(self):
        return '\n\n'.join(a + '\n\n' + repr(eval(a))
                        for a in self.args)
```

Planets Data

行星数据

Here we will use the Planets dataset, available via the [Seaborn package](#) (see [Visualization With Seaborn](#)). It gives information on planets that astronomers have discovered around other stars (known as extrasolar planets or exoplanets for short). It can be downloaded with a simple Seaborn command:

这里我们会使用Seaborn提供的行星数据 (参见[使用Seaborn进行可视化](#))。这个数据集提供了天文学家发现的其他恒星的行星的数据 (被称为太阳系外行星)。数据就可以简单的使用一个Seaborn命令来下载:

```
In [2]: import seaborn as sns
planets = sns.load_dataset('planets')
planets.shape
```

Out[2]: (1035, 6)

```
In [3]: planets.head()
```

```
Out[3]:
```

	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	110.82	2007
4	Radial Velocity	1	516.220	10.50	119.47	2009

This has some details on the 1,000+ extrasolar planets discovered up to 2014.

直到2014年已经有超过1000个太阳系外行星的数据。

Simple Aggregation in Pandas

在Pandas中进行简单聚合

Earlier, we explored some of the data aggregations available for NumPy arrays ([Aggregations: Min, Max, and Everything in Between](#)). As with a one-dimensional NumPy array, for a Pandas `Series`, the aggregates return a single value:

上章中, 我们已经介绍了NumPy数组的数据聚合操作 ([聚合: Min, Max, 以及其他](#))。正如一维NumPy数组, Pandas的 `Series` 的聚合结果是一个标量:

```
In [4]: rng = np.random.RandomState(42)
ser = pd.Series(rng.rand(5))
ser
```

Out[4]:

```
0    0.374540
1    0.959714
2    0.731994
3    0.598658
4    0.156019
dtype: float64
```

```
In [5]: ser.sum()
```

```
Out[5]: 2.811925491708157
```

```
In [6]: ser.mean()
```

```
Out[6]: 0.562385983416314
```

For a `DataFrame`, by default the aggregates return results within each column:

对于 `DataFrame` 来说, 默认情况下是每个列进行聚合的结果:

```
In [7]: df = pd.DataFrame({'A': rng.rand(5),
                        'B': rng.rand(5)})
df
```

Out[7]:

	A	B
0	0.155995	0.020584
1	0.058084	0.96910
2	0.866176	0.323443
3	0.601115	0.212339
4	0.708073	0.181825

```
In [8]: df.mean()
```

```
Out[8]: A    0.477888
        B    0.443420
dtype: float64
```

By specifying the `axis` parameter, you can instead aggregate within each row:

通过指定 `axis` 参数, 可以为每一行进行聚合操作:

```
In [9]: df.mean(axis='columns')
```

Out[9]:

```
0    0.989299
1    0.513997
2    0.849389
3    0.485727
4    0.444949
dtype: float64
```

Pandas `Series` and `DataFrame` s include all of the common aggregates mentioned in [Aggregations: Min, Max, and Everything in Between](#); in addition, there is a convenience method `describe()` that computes several common aggregates for each column and returns the result. Let's use this on the Planets data, for now dropping rows with missing values:

Pandas的 `Series` 和 `DataFrame` 包括了所有我们在[聚合: Min, Max, 以及其他](#)中介绍过的通用聚合操作; 而且Pandas还提供了很方便的 `describe()` 可以用来对每个列计算这些通用的聚合结果。让我们在我们的行星数据集上使用这个函数, 暂时先移除含有空值的行:

```
In [10]: planets.dropna().describe()
```

```
Out[10]:
```

	number	orbital period	mass	distance	year
count	498.000000	498.000000	498.000000	498.000000	498.000000
mean	1.73494	835.778671	2.509320	52.068213	2007.377510
std	1.17572	1469.128259	3.636274	46.596041	4.167284
min	1.00000	1.328300	0.003600	1.350000	1989.000000
25%	1.00000	38.272250	0.212500	24.497500	2005.000000
50%	1.00000	357.000000	1.245000	39.940000	2009.000000
75%	2.00000	999.600000	2.867500	69.332500	2011.000000
max	6.00000	17327.500000	25.000000	354.000000	2014.000000

This can be a useful way to begin understanding the overall properties of a dataset. For example, we see in the 'year' column that although exoplanets were discovered as far back as 1989, half of all known exoplanets were not discovered until 2010 or after. This is largely thanks to the *Kepler* mission, which is a space-based telescope specifically designed for finding eclipsing planets around other stars.

对于开始理解数据集的整体情况来说, 这是一个非常有用的方法。例如, 在发现年份 'year' 列上, 结果显示, 虽然第一颗太阳系外行星是1989年发现的, 但是一半的行星直到2010年以后才被发现。这多亏了开普勒Kepler计划, 它是一个太空望远镜, 专门设计用来寻找其他恒星的周围运行行星的。

The following table summarizes some other built-in Pandas aggregations:

聚合函数	描述
<code>count()</code>	元素个数
<code>first()</code> , <code>last()</code>	第一个和最后一个元素
<code>mean()</code> , <code>median()</code>	平均值和中位数
<code>min()</code> , <code>max()</code>	最小值和最大值
<code>std()</code> , <code>var()</code>	标准差和方差
<code>mad()</code>	平均绝对偏差
<code>prod()</code>	所有元素的乘积
<code>sum()</code>	所有元素的总和

These are all methods of `DataFrame` and `Series` objects.

它们都是 `DataFrame` 和 `Series` 对象的方法。

To go deeper into the data, however, simple aggregates are often not enough. The next level of data summarization is the `groupby` operation, which allows you to quickly and efficiently compute aggregates on subsets of data.

然而要深入了解数据, 简单的聚合经常是不够的。 `groupby` 操作为我们提供更高层次的概括功能, 通过它能很快和有效地计算子数据集的聚合数据。

GroupBy: Split, Apply, Combine

分组: 拆分、应用、组合

Simple aggregations can give you a flavor of your dataset, but often we need to aggregate conditionally on some label or index: this is implemented in the so-called "groupby" operation. The name "groupby" comes from a command in the SQL database language, but it perhaps more illuminative to think of it in the terms first coined by Hadley Wickham of Rstats fame: *split*, *apply*, *combine*.

简单的聚合可以提供数据集的基础特征, 但是通常我们希望依据一些标签或索引条件进行聚合操作; 这可以通过 `groupby` 操作实现。 "group by"的名称来自于SQL, 但是将它想成是由Hadley Wickham首先创造的R数据统计术语会更加合适: *拆分*、*应用*、*组合*。

Split, apply, combine

拆分、应用、组合

A canonical example of this split-apply-combine operation, where the "apply" is a summation aggregation, is illustrated in this figure:

作为一个具体例子, 我们来看一下使用Pandas来实现上面的这些计算, 首先创建一个输入 `DataFrame` :

```
In [11]: df = pd.DataFrame({'key': ['A', 'B', 'C', 'A', 'B', 'C'],
                        'data': ['A', 'B', 'C', 'A', 'B', 'C']})
df
```

Out[11]:

	key	data
0	A	A
1	B	1
2	C	2
3	A	3
4	B	4
5	C	5

The most basic split-apply-combine operation can be computed with the `groupby()` method of `DataFrame` s, passing the name of the desired key column:

最基础的拆分-应用-组合操作可以使用 `DataFrame` 的 `groupby()` 方法来实现, 方法中传递作为键来运算的列名:

```
In [12]: df.groupby('key')
```

```
Out[12]: <pandas.core.groupby.generic.DataFrameGroupBy object at 0x7fd196fa035e>
```

Notice that what is returned is not a set of `DataFrame` s, but a `DataFrameGroupBy` object. This object is where the magic is; you can think of it as a special view of the `DataFrame`, which is poised to dig into the groups but does no actual computation until the aggregation is applied. This "lazy evaluation" approach means that common aggregates can be implemented very efficiently in a way that is almost transparent to the user.

上面运行的结果并不是一个 `DataFrame`, 而是一个 `DataFrameGroupBy` 对象。这个对象就是上述步骤魔力的所在: 你可以认为它是 `DataFrame` 对象的一个特殊的视图, 使用它可以很容易的检索分组的数据, 但是除非聚合操作发生, 否则它不会进行真实的运算。这种"懒运算"的方式意味着通常的聚合可以实现得非常的高效, 而对于用户来说几乎是透明的。

To produce a result, we can apply an aggregate to this `DataFrameGroupBy` object, which will perform the appropriate `apply/combine` steps to produce the desired result:

要产生结果, 我们可以将一个聚合操作应用到该 `DataFrameGroupBy` 对象上, 这样就会在后台执行应用/组合的步骤, 并产生需要的结果:

```
In [13]: df.groupby('key').sum()
```

```
Out[13]:
```

	key	data
	A	3
	B	5
	C	7

The `sum()` method is just one possibility here; you can apply virtually any common Pandas or NumPy aggregation function, as well as virtually any valid `DataFrame` operation, as we will see in the following discussion.

`sum()` 方法仅仅是其中一个可能的操作; 你可以在这里应用几乎所有的Pandas或NumPy的通用聚合函数, 也可以应用集合所有正确的 `DataFrame` 操作, 我们将在下面马上就会看到。

The GroupBy object

GroupBy 对象

The `GroupBy` object is a very flexible abstraction. In many ways, you can simply treat it as if it's a collection of `DataFrame` s, and it does the difficult things under the hood. Let's see some examples using the Planets data.

`GroupBy` 对象是一个很灵活的抽象。在很多情况下, 你可以将它简单的看成 `DataFrame` 的集合, 它在底层做了很多复杂的工作。我们使用行星数据集来看几个例子。

Perhaps the most important operations made available by a `GroupBy` are *aggregate*, *filter*, *transform*, and *apply*. We'll discuss each of these more fully in [Aggregate, Filter, Transform, Apply](#), but before that let's introduce some of the other functionality that can be used with the basic `GroupBy` operation.

也许对于 `GroupBy` 对象最重要的操作是聚合、过滤、转换和应用。我们会在[聚合、过滤、转换、应用](#)中逐个介绍它们, 在这之前首先介绍一些其他用于 `GroupBy` 对象的基础操作。

Column indexing

The `GroupBy` object supports column indexing in the same way as the `DataFrame`, and returns a modified `GroupBy` object. For example:

`GroupBy` 对象支持列索引, 与 `DataFrame` 相同, 返回的是修改后的 `GroupBy` 对象。例如:

```
In [14]: planets.groupby('method')
```

```
Out[14]: <pandas.core.groupby.generic.DataFrameGroupBy object at 0x7fd196fa035e>
```

```
In [15]: planets.groupby('method')['orbital_period']
```

```
Out[15]: <pandas.core.groupby.generic.SeriesGroupBy object at 0x7fd19708dddb>
```

Here we've selected a particular `Series` group from the original `DataFrame` group by reference to its column name. As with the `GroupBy` object, no computation is done until we call some aggregate on the object:

上例中我们在原始的 `DataFrame` 中选择了特定的 `Series`。这个 `Series` 是按照提供的列名进行分组的。当然, `GroupBy` 对象在调用聚合操作之前是不会进行计算的:

```
In [16]: planets.groupby('method')['orbital_period'].median()
```

```
Out[16]: method
Astrometry                631.180000
Eclipse Timing Variations 4243.500000
Imaging                   27590.000000
Microlensing              3360.000000
Orbital Brightness Modulation 9.342807
Pulsar Timing             66.543900
Pulsation Timing Variations 1170.000000
Radial Velocity           308.200000
Transit                   5.734932
Transit Timing Variations 57.011000
Name: orbital_period, dtype: float64
```

This gives an idea of the general scale of orbital periods (in days) that each method is sensitive to.

结果给出了一个不同测量方法对公转周期进行测量的大概范围。

Iteration over groups

The `GroupBy` object supports direct iteration over the groups, returning each group as a `Series` or `DataFrame`:

`GroupBy` 对象支持在分组上直接进行迭代, 每次迭代返回分组的 `Series` 或 `DataFrame` 对象:

```
In [17]: for (method, group) in planets.groupby('method'):
        print("{} ({}s) shape={}".format(method, group.shape))
```

This can be useful for doing certain things manually, though it is often much faster to use the built-in `apply` functionality, which we will discuss momentarily.

这种做法在某些需要手动实现的情况下很有用, 虽然通常来说使用内置的 `apply` 函数会很快很多, 我们马上会介绍到 `apply` 函数。

Dispatch methods

扩展方法

Through some Python class magic, any method not explicitly implemented by the `GroupBy` object will be passed through and called on the groups, whether they are `DataFrame` or `Series` objects. For example, you can use the `describe()` method of `DataFrame` s to perform a set of aggregations that describe each group in the data:

通过一些Python面向对象的技术技巧, 任何非显式定义在 `GroupBy` 对象上的方法, 无论是 `DataFrame` 还是 `Series` 对象的, 都可以给分组来调用。例如, 你可以对数据分组上调用 `DataFrame` 的 `describe()` 方法, 对所有非聚合运算的聚合会运算:

译者注: 作者下面代码多加了 `unstack()` 方法, 应该是笔误。

```
In [18]: planets.groupby('method')['year'].describe()
```

```
Out[18]:
```

	count	mean	std	min	25%	50%	75%	max	
	Astrometry	2.0	2011.500000	2.121320	2010.0	2010.75	2011.5	2012.25	2013.0
	Eclipse Timing Variations	9.0	2010.000000	1.414214	2008.0	2009.00	2010.0	2011.00	2012.0
	Imaging	38.0	2009.131579	2.781901	2004.0	2008.00	2009.0	2010.00	2013.0
	Microlensing	23.0	2009.782609	2.859697	2004.0	2008.00	2010.0	2012.00	2013.0
	Orbital Brightness Modulation	3.0	2011.666667	1.154701	2011.0	2011.00	2011.0	2012.00	2013.0
	Pulsar Timing	5.0	1998.400000	0.384510	1992.0	1992.00	1994.0	2003.00	2011.0
	Pulsation Timing Variations	1.0	2007.000000	NaN	2007.0	2007.00	2007.0	2007.00	2007.0
	Radial Velocity	553.0	2007.518897	4.249052	1989.0	2005.00	2009.0	2011.00	2014.0
	Transit	397.0	2010.236716	2.077867	2002.0	2010.00	2012.0	2013.00	2014.0
	Transit Timing Variations	4.0	2012.500000	1.290994	2011.0	2011.75	2012.5	2013.25	2014.0

Looking at this table helps us to better understand the data: for example, the vast majority of planets have been discovered by the Radial Velocity and Transit methods, though the latter only became common (due to new, more accurate telescopes) in the last decade. The newest methods seem to be Transit Timing Variation and Orbital Brightness Modulation, which were not used to discover a new planet until 2011.

查看上表, 能帮助我们更好的理解数据: 例如, 发现行星最多的方法是径向速度和凌日法, 虽然后者是近十年才变得普遍 (因为新的更精准的望远镜的使用)。最新的方法应该是凌日时间变分和轨道亮度调制法, 它们直至2011年才开始发现新的行星。

This is just one example of the utility of dispatch methods. Notice that they are applied to each individual group, and the results are then combined within `GroupBy` and returned. Again, any valid `DataFrame` / `Series` method can be used on the corresponding `GroupBy` object, which allows for a very flexible and powerful operations!

这只是一个使用扩展方法的例子。你需要知道的是这些方法会被应用到每一个独立的分组上, 然后计算得到的结果会在 `GroupBy` 对象中合并并返回。再次提示, 任何正确的 `DataFrame` 或 `Series` 方法都能在相应的 `GroupBy` 对象上使用, 这种扩展方法的方式提供了非常灵活及强大的操作。

Aggregate, filter, transform, apply

聚合、过滤、转换、应用

The preceding discussion focused on aggregation for the combine operation, but there are more options available. In particular, `GroupBy` objects have `aggregate()`, `filter()`, `transform()`, and `apply()` methods that efficiently implement a variety of useful operations before combining the group data.

前面的讨论聚焦在组合操作相应的聚合函数上, 但实际上还有更多的可能选项。特别是 `GroupBy` 对象有 `aggregate()`、`filter()`、`transform()` 和 `apply()` 方法, 它们能在组合分组数据之前有效地实现大量有用的操作。

For the purpose of the following subsections, we'll use this `DataFrame`:

对于下面的部分内容, 我们将使用下述的 `DataFrame`:

```
In [19]: rng = np.random.RandomState(0)
df = pd.DataFrame({'key': ['A', 'B', 'C', 'A', 'B', 'C'],
                  'data1': range(6),
                  'data2': rng.randint(0, 10, 6)},
                  columns = ['key', 'data1', 'data2'])
df
```

Out[19]:

	key	data1	data2
0	A	0	5
1	B	1	1
2	C	2	3
3	A	3	3
4	B	4	7
5	C	5	9

Aggregation

We're now familiar with `GroupBy` aggregations with `sum()`, `median()`, and the like, but the `aggregate()` method allows for even more flexibility. It can take a string, a function, or a list thereof, and compute all the aggregates at once. Here is a quick example combining all these:

我们已然熟悉 `GroupBy` 的 `sum()`、`median()` 等方法计算聚合的做法, 但是 `aggregate()` 方法能提供更多的灵活性。它能接受字符串、函数或者一个列表, 然后一次性计算所有的聚合结果。下面是一个简单的例子:

```
In [20]: df.groupby('key').aggregate(['min', np.median, max])
```

Out[20]:

	data1	data2				
	min	median	max	min	median	max
A	0	1.5	3	0	4.0	5
B	1	2.5	4	0	3.5	7
C	2	3.5	5	3	6.0	9

Another useful pattern is to pass a dictionary mapping column names to operations to be applied on that column:

还可以将一个字典, 里面是列名与操作的对对应关系, 传递给 `aggregate()` 来进行一次性的聚合运算:

```
In [21]: df.groupby('key').aggregate({'data1': 'min',
                                    'data2': 'max'})
```

Out[21]:

	data1	data2
	key	
A	0	5
B	1	7
C	2	9

Filtering

A filtering operation allows you to drop data based on the group properties. For example, we might want to keep all groups in which the standard deviation is larger than some critical value:

过滤操作能让你在分组数据上移除你不需要的数据。例如, 我们可能希望保留标准差大于某个阈值的所有的分组:

译者注: 你可以认为 `filter()` 类似于SQL中的HAVING。

```
In [22]: def filter_func(x):
        return x['data2'].std() > 4
display('df', "df.groupby('key').std()", "df.groupby('key').filter(filter_func)")
```

Out[22]:

	key	data1	data2
0	A	0	5
1	B	1	0
2	C	2	3
3	A	3	3
4	B	4	7
5	C	5	9

The filter function should return a Boolean value specifying whether the group passes the filtering. Here because group A does not have a standard deviation greater than 4, it is dropped from the result.

用来进行过滤的函数必须返回一个布尔值, 表示分组是否能够通过过滤条件。上例中A分组的标准差不是大于4, 因此整个分组在结果中被移除了。

Transformation

转换

While aggregation must return a reduced version of

In [26]: display('df', "df.groupby(df['key']).sum()")

Out[26]:

```
df
df.groupby(df['key']).sum()
```

	key	data1	data2		data1	data2
0	A	0	5	key		
1	B	1	0	A	3	8
2	C	2	3	B	5	7
3	A	3	3	C	7	12
4	B	4	7			
5	C	5	9			

A dictionary or series mapping index to group

使用字典或映射索引的序列来分组

Another method is to provide a dictionary that maps index values to the group keys:

还有一种方法是提供一个字典，将索引值映射成分组键：

In [27]: df2 = df.set_index('key')
mapping = {'A': 'vowel', 'B': 'consonant', 'C': 'consonant'}
display('df2', 'df2.groupby(mapping).sum()')

Out[27]:

```
df2
df2.groupby(mapping).sum()
```

	data1	data2		data1	data2
key			consonant	12	19
A	0	5	vowel	3	8
B	1	0			
C	2	3			
A	3	3			
B	4	7			
C	5	9			

Any Python function

任何Python函数

Similar to mapping, you can pass any Python function that will input the index value and output the group:

类似映射，你可以传递任何python函数将输入的索引值变成输出的分组键：

In [28]: display('df2', 'df2.groupby(str.lower).mean()')

Out[28]:

```
df2
df2.groupby(str.lower).mean()
```

	data1	data2		data1	data2
key			a	1.5	4.0
A	0	5	b	2.5	3.5
B	1	0	c	3.5	6.0
C	2	3			
A	3	3			
B	4	7			
C	5	9			

A list of valid keys

正确的列表

Further, any of the preceding key choices can be combined to group on a multi-index:

还有，任何前面的多个分组键可以组合并输出成一个多重索引的结果：

In [29]: df2.groupby([str.lower, mapping]).mean()

Out[29]:

	data1	data2
a	vowel	1.5 4.0
b	consonant	2.5 3.5
c	consonant	3.5 6.0

Grouping example

分组例子

As an example of this, in a couple lines of Python code we can put all these together and count discovered planets by method and by decade:

作为分组的例子，我们将前面介绍的内容用几行Python代码写出来用于计算通过不同方法在不同年代发现的行星的个数：

In [30]: decade = 10 * (planets['year'] // 10)
decade = decade.astype(str) + 's'
decade.name = 'decade'
planets.groupby(['method', decade])['number'].sum().unstack().fillna(0)

Out[30]:

decade		1980s	1990s	2000s	2010s
	method				
	Astrometry	0.0	0.0	0.0	2.0
	Eclipse Timing Variations	0.0	0.0	5.0	10.0
	Imaging	0.0	0.0	29.0	21.0
	Microlensing	0.0	0.0	12.0	15.0
	Orbital Brightness Modulation	0.0	0.0	0.0	5.0
	Pulsar Timing	0.0	9.0	1.0	1.0
	Pulsation Timing Variations	0.0	0.0	1.0	0.0
	Radial Velocity	1.0	52.0	475.0	424.0
	Transit	0.0	0.0	64.0	712.0
	Transit Timing Variations	0.0	0.0	0.0	9.0

This shows the power of combining many of the operations we've discussed up to this point when looking at realistic datasets. We immediately gain a coarse understanding of when and how planets have been discovered over the past several decades!

这个例子展示了我们结合前面介绍过的多种操作之后，我们能在真实的数据集上完成多强大的操作。我们立即获得了过去几十年间我们是如何发现行星的大概统计。

Here I would suggest digging into these few lines of code, and evaluating the individual steps to make sure you understand exactly what they are doing to the result. It's certainly a somewhat complicated example, but understanding these pieces will give you the means to similarly explore your own data.

作者建议你深入研究上面的几行代码，逐步的执行它们，直到你完全理解了这些代码是如何最终产生结果的。当然上面是一个稍微复杂的例子，但理解这个例子会让你在研究自己的数据集时知道如何进行操作。

< [组合数据集: Merge 和 Join](#) | [目录](#) | [数据透视表](#) >

 [Open in Colab](#)