

- 3.6 Pandas层级索引
 - 3.6.1 多级索引Series
 - 3.6.1.1 笨办法
 - 3.6.1.2 好办法：Pandas多级索引
 - 3.6.1.3 高维数据的多级索引(重要)
 - 3.6.2 多级索引的创建方法
 - 3.6.2.1 显式地创建多级索引
 - 3.6.2.2 多级索引的等级名称
 - 3.6.2.3 多级列索引
 - 3.6.3 多级索引的取值与切片
 - 3.6.3.1 Series多级索引
 - 3.6.3.2 DataFrame多级索引
 - 3.6.4 多级索引行列转换
 - 3.6.4.1 有序的索引和无序的索引
 - 3.6.4.2 索引stack与unstack
 - 3.6.4.3 索引的设置与重置
 - 3.6.5 多级索引的数据累计方法
 - 3.6.6 Panel数据



3.6 Pandas层级索引

《Python数据分析》课程讲义pandas模块内容，参考教材为《Python数据科学手册》
<https://www.jianshu.com/p/760cd4f46c8d>

当目前为止，接触的都是一维数据和二维数据，用 Pandas 的 `Series` 和 `DataFrame` 对象就可以存储。但也经常会遇到存储多维数据的需求，**数据索引超过一两个键**。为此，Pandas 提供了 `Panel` 和 `Panel4D` 对象解决三维数据与四维数据(详情请参见 3.7 节)。而在实践中，更直观的形式是通过**层级索引**(hierarchical indexing, 也被称为**多级索引**, multi-indexing)配合多个有不同**等级**(level)的一级索引一起使用，这样就可以将高维数组转换成类似一维 `Series` 和二维 `DataFrame` 对象的形式。

多级索引：在一个轴axis上有多个(两个以上)的索引，能够以低维度形式来表示高维度的数据。单级索引是Index对象，多级索引是MultiIndex对象。

在这一节中，将介绍创建 `MultiIndex` 对象的方法，多级索引数据的取值、切片和统计值的计算，以及普通索引与层级索引的转换方法。首先导入 Pandas 和 NumPy：

```
[1]:
```

```
1 import pandas as pd
2 import numpy as np
```

3.6.1 多级索引Series

让看看如何用一维的 `Series` 对象表示二维数据—用一系列包含特征与数值的数据点来简单演示。

3.6.1.1 笨办法

假设想要分析美国各州在两个不同年份的数据。如果用前面介绍的 Pandas 工具来处理，那么可能会用一个 Python 元组来表示索引：

[3]:

```
1 index = [('California', 2000), ('California', 2010),
2          ('New York', 2000), ('New York', 2010),
3          ('Texas', 2000), ('Texas', 2010)]
4 populations = [33871648, 37253956,
5                18976457, 19378102,
6                20851820, 25145561]
7 pop = pd.Series(populations, index=index)
8 print(pop)
9
10 (California, 2000)    33871648
11 (California, 2010)    37253956
12 (New York, 2000)     18976457
13 (New York, 2010)     19378102
14 (Texas, 2000)        20851820
15 (Texas, 2010)        25145561
16 dtype: int64
```

通过**元组构成的多级索引**，可以直接在 `Series` 上取值或用切片查询数据：

[4]:

```
1 print(pop[('California', 2000):('Texas', 2000)])
2
3 (California, 2000)    33871648
4 (California, 2010)    37253956
5 (New York, 2000)     18976457
6 (New York, 2010)     19378102
7 (Texas, 2000)        20851820
8 dtype: int64
```

但是这么做很不方便。假如想要选择所有 2000 年的数据，那么就得用一些比较复杂的(可能也比较慢的)清理方法了：

[5]:

```

1 print(pop[[i for i in pop.index if i[1] == 2000]])
2
3 (California, 2000)    33871648
4 (New York, 2000)     18976457
5 (Texas, 2000)        20851820
6 dtype: int64

```

这么做虽然也能得到需要的结果，但是与 Pandas 令人爱不释手的切片语法相比，这种方法确实不够简洁（在处理较大的数据时也不够高效）。

3.6.1.2 好办法：Pandas多级索引

好在 Pandas 提供了更好的解决方案。用元组表示索引其实是多级索引的基础，Pandas 的 `MultiIndex` 类型提供了更丰富的操作方法。可以用元组创建一个多级索引，如下所示：

[7]:

```

1 print(pop)
2 print('-----pd.MultiIndex.from_tuples(index)-----')
3 index = pd.MultiIndex.from_tuples(index)
4 print(index)
5
6 (California, 2000)    33871648
7 (California, 2010)    37253956
8 (New York, 2000)     18976457
9 (New York, 2010)     19378102
10 (Texas, 2000)        20851820
11 (Texas, 2010)        25145561
12 dtype: int64
13 -----pd.MultiIndex.from_tuples(index)-----
14 MultiIndex([('California', 2000),
15             ('California', 2010),
16             ('New York', 2000),
17             ('New York', 2010),
18             ('Texas', 2000),
19             ('Texas', 2010)],
20            )

```

会发现 `MultiIndex` 里面有一个 `levels` 属性表示索引的**等级**（索引的维度）——这样做可以将州名和年份作为每个数据点的不同**标签**。

总维度= 数据(列)维度 + 索引(行)维度。比如上述例子，数据维度为1，索引维度为2，那么总维度为3

如果将前面创建的 `pop` 的索引重置(`reindex`)为 `MultiIndex`，就会看到层级索引：

[10]:

```

1 pop = pop.reindex(index)
2 print(pop)
3
4 California 2000    33871648
5             2010    37253956
6 New York   2000    18976457
7             2010    19378102
8 Texas      2000    20851820
9             2010    25145561
10 dtype: int64

```

其中前两列表示 `Series` 的多级索引值(二维)，第三列是数据(一维)。会发现有些行仿佛缺失了第一列数据——这其实是多级索引的表现形式，每个空格与上面的索引相同。现在可以直接用第二个索引获取 2010 年的全部数据，与 Pandas 的多维切片 查询用法一致：

[11]:

```

1 print(pop[:, 2010])
2 #print(pop['Texas'])
3
4 California    37253956
5 New York      19378102
6 Texas         25145561
7 dtype: int64

```

结果是单索引的数组，正是需要的。与之前的元组索引相比，多级索引的语法更简洁。（操作更方便！）下面继续介绍层级索引的取值操作方法。

3.6.1.3 高维数据的多级索引(重要)

可能已经注意到，其实完全可以用一个带行列索引的简单 `DataFrame` 代替前面的多级索引。其实 Pandas 已经实现了类似的功能。`unstack()` 方法可以快速将一个多级索引的 `Series` 转化为普通索引的 `DataFrame`。

1. 行列索引转化，`.unstack()` 数据升维，索引降维(最内层索引)。而 `.stack()` 数据降维，索引升维(最内层索引)。
2. 每列的元素总个数 = 各维度行索引数的乘积。 $6 = 2 * 3$
3. 元素总个数 = 各维度行索引数*列索引数。 $6 = 2 * 3 * 1$

[13]:

```

1 print(pop)
2 print('---pop.unstack()---')
3 pop_df = pop.unstack()
4 print(pop_df)
5
6 California 2000    33871648
7             2010    37253956
8 New York   2000    18976457
9             2010    19378102
10 Texas      2000    20851820
11             2010    25145561

```

```

12 dtype: int64
13 ---pop.unstack()----
14           2000      2010
15 California 33871648 37253956
16 New York   18976457 19378102
17 Texas      20851820 25145561

```

当然了，也有 `stack()` 方法实现相反的效果：

[14]:

```

1 print(pop_df)
2 print('---pop.stack()----')
3 print(pop_df.stack())
4
5           2000      2010
6 California 33871648 37253956
7 New York   18976457 19378102
8 Texas      20851820 25145561
9 ---pop.stack()----
10 California 2000      33871648
11           2010      37253956
12 New York   2000      18976457
13           2010      19378102
14 Texas      2000      20851820
15           2010      25145561
16 dtype: int64

```

可能会纠结于为什么要花时间研究层级索引。其实理由很简单：如果可以用含多级索引的一维 `Series` 数据表示二维数据，那么就可以用 `Series` 或 `DataFrame` 表示三维甚至更高维度的数据。多级索引每增加一级，就表示数据增加一维，利用这一特点就可以轻松表示任意维度的数据了。假如要增加一列显示每年各州的人口统计指标(例如 18 岁以下的人口)，那么对于这种带有 `MultiIndex` 的对象，增加一列就像 `DataFrame` 的操作一样简单：

[15]:

```

1 pop_df = pd.DataFrame({'total': pop,
2                         'under18': [9267089, 9284094,
3                                     4687374, 4318033,
4                                     5906301, 6879014]})
5 print(pop_df)
6 print('---pop.stack()----')
7 print(pop_df.stack())
8
9           total  under18
10 California 2000 33871648 9267089
11           2010 37253956 9284094
12 New York   2000 18976457 4687374
13           2010 19378102 4318033
14 Texas      2000 20851820 5906301
15           2010 25145561 6879014
16 ---pop.stack()----

```

```

17 California 2000 total 33871648
18             under18 9267089
19             2010 total 37253956
20             under18 9284094
21 New York   2000 total 18976457
22             under18 4687374
23             2010 total 19378102
24             under18 4318033
25 Texas      2000 total 20851820
26             under18 5906301
27             2010 total 25145561
28             under18 6879014
29 dtype: int64

```

另外，所有在 3.4 节介绍过的通用函数和其他功能也同样适用于层级索引。可以计算上面数据中 18 岁以下的人口占总人口的比例：

[16]:

```

1 f_u18 = pop_df['under18'] / pop_df['total']
2 print(f_u18)
3 print('-----f_u18.unstack()-----')
4 print(f_u18.unstack())
5
6 California 2000    0.273594
7             2010    0.249211
8 New York   2000    0.247010
9             2010    0.222831
10 Texas      2000    0.283251
11             2010    0.273568
12 dtype: float64
13 -----f_u18.unstack()-----
14             2000    2010
15 California 0.273594 0.249211
16 New York   0.247010 0.222831
17 Texas      0.283251 0.273568

```

同样，也可以快速浏览和操作高维数据。

3.6.2 多级索引的创建方法

为 `Series` 或 `DataFrame` 创建多级索引最直接的办法就是将 `index` 参数设置为至少二维的索引数组，如下所示：

[83]:

```

1 df = pd.DataFrame(np.random.rand(4, 2), # 4*2 = 8
2                   index=[['English', 'English', 'Chinese', 'Chinese'],
3                           ['like', 'dislike', 'like', 'dislike']], # 4 + 4
4                   columns=['girl', 'boy'])
5 print(df)
6 print('---MultiIndex.from_product()笛卡尔积外积---')

```

```

7 df2 = pd.DataFrame(np.random.rand(4, 2),
8                     index=pd.MultiIndex.from_product(['English', 'Chinese'],
9                                                       ['like', 'dislike'])), # 2*2
10                     columns= ['girl', 'boy']) # 2
11 print(df2) # 创建多级 行 索引
12
13      girl    boy
14 English like  0.705040  0.611248
15      dislike  0.661523  0.745367
16 Chinese like  0.749331  0.339376
17      dislike  0.124760  0.102350
18 ---MultiIndex.from_product()---
19      girl    boy
20 English like  0.247010  0.956254
21      dislike  0.468136  0.120788
22 Chinese like  0.188142  0.398446
23      dislike  0.229750  0.071471

```

`MultiIndex` 的创建工作将在后台完成。同理，如果把将元组作为键的字典传递给 Pandas，Pandas 也会默认转换为 `MultiIndex`：

[84]:

```

1 data = {('California', 2000): 33871648,
2         ('California', 2010): 37253956,
3         ('Texas', 2000): 20851820,
4         ('Texas', 2010): 25145561,
5         ('New York', 2000): 18976457,
6         ('New York', 2010): 19378102}
7 print(pd.Series(data))
8
9 California 2000    33871648
10           2010    37253956
11 Texas      2000    20851820
12           2010    25145561
13 New York   2000    18976457
14           2010    19378102
15 dtype: int64

```

但是有时候显式地创建 `MultiIndex` 也是很有用的，下面来介绍一些创建方法。

3.6.2.1 显式地创建多级索引

可以用 `pd.MultiIndex` 中的类方法更加灵活地构建多级索引。下面介绍4种常见的方式：

1. 通过一个有不同等级的若干简单数组组成的列表来构建 `MultiIndex`：

[85]:

```

1 print(pd.MultiIndex.from_arrays([ ['a', 'a', 'b', 'b'], [1, 2, 1, 2] ]))
2
3 MultiIndex(levels=[['a', 'b'], [1, 2]],
4             codes=[[0, 0, 1, 1], [0, 1, 0, 1]])

```

1. 通过包含多个索引值的元组构成的列表创建 `MultiIndex` :

[86]:

```

1 print(pd.MultiIndex.from_tuples([('a', 1), ('a', 2), ('b', 1), ('b', 2)]))
2
3 MultiIndex([('a', 1),
4            ('a', 2),
5            ('b', 1),
6            ('b', 2)],
7            )

```

1. 通过两个索引的笛卡尔积(Cartesian product)创建 `MultiIndex` :

[87]:

```

1 print(pd.MultiIndex.from_product([['a', 'b'], [1, 2]]))
2
3 MultiIndex([('a', 1),
4            ('a', 2),
5            ('b', 1),
6            ('b', 2)],
7            )

```

1. 提供 `levels` (包含每个等级的索引值列表的列表)和 `labels` (包含每个索引值标签列表的列表)创建 `MultiIndex` :

[88]:

```

1 print(pd.MultiIndex(levels=[['a', 'b'], [1, 2]],
2                       codes=[[0, 0, 1, 1], [0, 1, 0, 1]]))
3
4 MultiIndex([('a', 1),
5            ('a', 2),
6            ('b', 1),
7            ('b', 2)],
8            )

```

在创建 `Series` 或 `DataFrame` 时, 可以将这些对象作为 `index` 参数, 或者通过 `reindex` 方法更新 `Series` 或 `DataFrame` 的索引.

3.6.2.2 多级索引的等级名称

给 `MultiIndex` 的等级加上名称会为一些操作提供便利. 可以在前面任何一个 `MultiIndex` 构造器中通过 `names` 参数设置等级名称, 也可以在创建之后通过索引的 `names` 属性来修改名称:

[89]:

```


```



```

1 print(pop)
2 print('-----pop.index.names-----')
3 pop.index.names = ['state', 'year']
4 print(pop)
5
6 California 2000    33871648
7             2010    37253956
8 New York   2000    18976457
9             2010    19378102
10 Texas     2000    20851820
11             2010    25145561
12 dtype: int64
13 -----pop.index.names-----
14 state      year
15 California 2000    33871648
16             2010    37253956
17 New York   2000    18976457
18             2010    19378102
19 Texas     2000    20851820
20             2010    25145561
21 dtype: int64

```

在处理复杂的数据时，为等级设置名称是管理多个索引值的好办法。

3.6.2.3 多级列索引

每个 `DataFrame` 的行与列都是对称的，也就是说既然有多级行索引，那么同样可以有多级列索引。让通过一份医学报告的模拟数据来演示：

[17]:

```

1 # 多级行列索引
2 index = pd.MultiIndex.from_product([[2013, 2014], [1, 2]],
3                                     names=['year', 'visit']) # 2*2
4 columns = pd.MultiIndex.from_product([['Bob', 'Guido', 'Sue'], ['HR', 'Temp']],
5                                     names=['subject', 'type']) # 3*2
6
7 print(index)
8 print(columns)
9
10 # 模拟数据
11 data = np.round(np.random.randn(4, 6), 1)
12 data[:, ::2] *= 10
13 data += 37
14 # 创建DataFrame
15 health_data = pd.DataFrame(data, index=index, columns=columns)
16 health_data
17
18 MultiIndex([(2013, 1),
19             (2013, 2),
20             (2014, 1),

```

```

21         (2014, 2)],
22         names=['year', 'visit'])
23 MultiIndex([( 'Bob',   'HR'),
24             ( 'Bob',   'Temp'),
25             ('Guido',   'HR'),
26             ('Guido',   'Temp'),
27             ( 'Sue',   'HR'),
28             ( 'Sue',   'Temp')]),
29         names=['subject', 'type'])

```

[17]:

	subject	Bob		Guido		Sue	
	type	HR	Temp	HR	Temp	HR	Temp
year	visit						
2013	1	47.0	35.7	24.0	37.3	35.0	35.9
	2	20.0	34.7	49.0	37.5	29.0	36.7
2014	1	41.0	38.3	34.0	36.0	37.0	37.3
	2	42.0	36.3	41.0	36.1	29.0	35.6

多级行列索引的创建非常简单。上面创建了一个简易的四维数据，四个维度分别为被检查人的姓名、检查项目、检查年份和检查次数。可以在列索引的第一级查询姓名，从而获取包含一个人(例如 Guido)全部检查信息的 **DataFrame**：

[91]:

```

1 health_data['Guido']

```

[91]:

	type	HR	Temp
year	visit		
2013	1	24.0	37.3
	2	49.0	37.5
2014	1	34.0	36.0
	2	41.0	36.1

如果想获取包含多种标签的数据，需要通过对多个维度(姓名、国家、城市等标签)的多次查询才能实现，这时使用多级行列索引进行查询会非常方便。

3.6.3 多级索引的取值与切片

对 `MultiIndex` 的取值和切片操作很直观，可以直接把索引看成额外增加的维度。先来介绍 `Series` 多级索引的取值与切片方法，再介绍 `DataFrame` 的用法。

3.6.3.1 Series多级索引

看看下面由各州历年人口数量创建的多级索引 `Series`：

[18]:

```
1 print(pop)
2
3 California 2000    33871648
4             2010    37253956
5 New York   2000    18976457
6             2010    19378102
7 Texas      2000    20851820
8             2010    25145561
9 dtype: int64
```

可以通过对多个级别索引值获取单个元素：

[93]:

```
1 print(pop['California', 2000])
2
3 33871648
```

`MultiIndex` 也支持**局部取值**(partial indexing)，即**只取索引的某一个层级，由外层到内层**。假如只取最高级的索引，那么会把内层的索引保留，获得的结果是一个新的 `Series`，未被选中的低层索引值会被保留：

[23]:

```
1 print(pop['Texas'])
2 print(pop[:, 2000])
3
4 2000    20851820
5 2010    25145561
6 dtype: int64
7 California 33871648
8 New York   18976457
9 Texas      20851820
10 dtype: int64
```

类似的还有局部切片，不过要求 `MultiIndex` 是按顺序排列的(就像将在 3.6.4 节介绍的那样)：

[24]:

```

1 print(pop.loc['California':'New York'])
2
3 California  2000      33871648
4              2010      37253956
5 New York    2000      18976457
6              2010      19378102
7 dtype: int64

```

如果索引已经排序，那么可以用**低层级的索引取值，高层索引需要使用:切片**：

[96]:

```

1 print(pop[:, 2000])
2
3 state
4 California      33871648
5 New York        18976457
6 Texas           20851820
7 dtype: int64

```

其他取值与数据选择的方法(详情请参见 3.3 节)也都起作用。下面的例子是通过布尔掩码选择数据：

[97]:

```

1 print(pop[pop > 22000000])
2
3 state      year
4 California  2000      33871648
5              2010      37253956
6 Texas       2010      25145561
7 dtype: int64

```

也可以用花哨的索引选择数据：

[98]:

```

1 print(pop[ ['California', 'Texas'] ])
2
3 state      year
4 California  2000      33871648
5              2010      37253956
6 Texas       2000      20851820
7              2010      25145561
8 dtype: int64

```

3.6.3.2 DataFrame多级索引

DataFrame 多级索引的用法与 **Series** 类似。还用之前的体检报告数据来演示：

[99]:

```

1 health_data

```

[99]:

	subject	Bob		Guido		Sue	
	type	HR	Temp	HR	Temp	HR	Temp
year	visit						
2013	1	47.0	35.7	24.0	37.3	35.0	35.9
	2	20.0	34.7	49.0	37.5	29.0	36.7
2014	1	41.0	38.3	34.0	36.0	37.0	37.3
	2	42.0	36.3	41.0	36.1	29.0	35.6

由于 `DataFrame` 的基本索引是列索引，因此 `Series` 中多级索引的用法到了 `DataFrame` 中就应用在列上了。例如，可以通过简单的操作获取 Guido 的心率数据：

[100]:

```
1 print(health_data['Guido', 'HR'])
2
3 year  visit
4 2013    1    24.0
5      2    49.0
6 2014    1    34.0
7      2    41.0
8 Name: (Guido, HR), dtype: float64
```

与单索引类似，在 3.3 节介绍的 `loc`、`iloc` 和 `ix` 索引器都可以使用，例如：

[101]:

```
1 health_data.iloc[:2, :2]
```

[101]:

	subject	Bob	
	type	HR	Temp
year	visit		
2013	1	47.0	35.7
	2	20.0	34.7

虽然这些索引器将多维数据当作二维数据处理，但是在 `loc` 和 `iloc` 中可以传递多个层级的索引元组，例如：

[102]:

```

1 # 行索引 和 列索引
2 print(health_data.loc[:, ('Bob', 'HR')])
3
4 year  visit
5 2013  1      47.0
6       2      20.0
7 2014  1      41.0
8       2      42.0
9 Name: (Bob, HR), dtype: float64

```

这种索引元组的用法不是很方便，如果在元组中使用切片还会导致语法错误：

[103]:

```

1 # health_data.loc[:, 1], (:, 'HR')]
2
3 SyntaxError: invalid syntax (Temp/ipykernel_17224/3311942670.py, line 1)
4   File "C:\Users\15152\AppData\Local\Temp/ipykernel_17224/3311942670.py", line 1
5     health_data.loc[:, 1], (:, 'HR')]
6                         ^
7 SyntaxError: invalid syntax

```

虽然可以用 Python 内置的 `slice()` 函数获取想要的切片，但是还有一种更好的办法，就是使用 `IndexSlice` 对象。Pandas 专门用它解决这类问题，例如：

[104]:

```

1 idx = pd.IndexSlice
2 health_data.loc[idx[:, 1], idx[:, 'HR']]

```

[104]:

	subject	Bob	Guido	Sue
	type	HR	HR	HR
year	visit			
2013	1	47.0	24.0	35.0
2014	1	41.0	34.0	37.0

和带多级索引的 `Series` 和 `DataFrame` 进行数据交互的方法有很多，但就像本书中的诸多工具一样，若想掌握，最好的办法就是使用它们！

- 基本索引 采用列索引方式；
- 索引器 采用行列索引方式； `health_data.iloc[:, :2]`
- 索引器 对行列多级索引会出错，索引引入 `IndexSlice` 对象。

3.6.4 多级索引行列转换

使用多级索引的关键是掌握有效数据转换的方法。Pandas 提供了许多操作，可以让数据在内容保持不变的同时，按照需要进行行列转换。之前用一个简短的例子演示过 `stack()` 和 `unstack()` 的用法，但其实还有许多合理控制层级行列索引的方法，让来一探究竟。

3.6.4.1 有序的索引和无序的索引

在前面的内容里，曾经简单提过多级索引排序，这里需要详细介绍一下。如果 `MultiIndex` 不是有序的索引，那么大多数切片操作都会失败。首先创建一个不按字典顺序(lexographically)排列的多级索引

`Series` :

[25]:

```
1 index = pd.MultiIndex.from_product([ ['a', 'c', 'b'], [1, 2] ])
2 data = pd.Series(np.random.rand(6), index=index)
3 data.index.names = ['char', 'int']
4 print(data)
5
6 char  int
7 a      1    0.073217
8        2    0.978929
9 c      1    0.594405
10       2    0.167412
11 b      1    0.786752
12       2    0.844419
13 dtype: float64
```

如果想对索引使用局部切片，那么错误就会出现：

[26]:

```
1 try:
2     data['a':'b']
3 except KeyError as e:
4     print(type(e))
5     print(e)
6
7 <class 'pandas.errors.UnsortedIndexError'>
8 'Key length (1) was greater than MultiIndex lexsort depth (0)'
```

尽管从错误信息里面看不出具体的细节，但问题是出在 `MultiIndex` 无序排列上。局部切片和许多其他相似的操作都要求 `MultiIndex` 的各级索引是有序的(即按照字典顺序由 A 至 Z)。为此，Pandas 提供了许多便捷的操作完成排序，如 `sort_index()` 和 `sortlevel()` 方法。用最简单的 `sort_index()` 方法来演示：

[28]:

```

1 data = data.sort_index()
2 print(data)
3
4 char  int
5 a      1      0.073217
6        2      0.978929
7 b      1      0.786752
8        2      0.844419
9 c      1      0.594405
10       2      0.167412
11 dtype: float64

```

索引排序之后，局部切片就可以正常使用了：

[29]:

```

1 print(data['a':'b'])
2
3 char  int
4 a      1      0.073217
5        2      0.978929
6 b      1      0.786752
7        2      0.844419
8 dtype: float64

```

3.6.4.2 索引stack与unstack

前文曾提过，可以将一个多级索引数据集转换成简单的二维形式，可以通过 `level` 参数设置转换的索引层级：

[109]:

```

1 print(pop)
2 print('-----unstack(level=0)---- ')
3 print(pop.unstack(level=0))
4
5 state      year
6 California 2000    33871648
7            2010    37253956
8 New York   2000    18976457
9            2010    19378102
10 Texas     2000    20851820
11           2010    25145561
12 dtype: int64
13 -----unstack(level=0)----
14 state California New York Texas
15 year
16 2000    33871648 18976457 20851820
17 2010    37253956 19378102 25145561

```

[110]:


```

1 print(pop)
2 print('-----unstack(level=1)---- ')
3 print(pop.unstack(level=1))
4
5 state      year
6 California 2000    33871648
7             2010    37253956
8 New York   2000    18976457
9             2010    19378102
10 Texas      2000    20851820
11             2010    25145561
12 dtype: int64
13 -----unstack(level=1)----
14 year          2000      2010
15 state
16 California 33871648 37253956
17 New York   18976457 19378102
18 Texas      20851820 25145561

```

`unstack()` 是 `stack()` 的逆操作，同时使用这两种方法让数据保持不变：

[111]:

```

1 print(pop.unstack().stack())
2
3 state      year
4 California 2000    33871648
5             2010    37253956
6 New York   2000    18976457
7             2010    19378102
8 Texas      2000    20851820
9             2010    25145561
10 dtype: int64

```

3.6.4.3 索引的设置与重置

层级数据维度转换的另一种方法是行列标签转换，可以通过 `reset_index` 方法实现。如果在上面的人口数据 `Series` 中使用该方法，则会生成一个列标签中包含之前索引标签 `state` 和 `year` 的 `DataFrame`。也可以用数据的 `name` 属性为列设置名称：

[30]:

```

1 print(pop)
2 print('-----pop.reset_index---- ')
3 pop_flat = pop.reset_index(name='population')
4 print(pop_flat)
5
6 California 2000    33871648
7             2010    37253956
8 New York   2000    18976457
9             2010    19378102
10 Texas      2000    20851820

```

```

11      2010      25145561
12 dtype: int64
13 -----pop.reset_index-----
14      level_0  level_1  population
15  0  California      2000      33871648
16  1  California      2010      37253956
17  2    New York      2000      18976457
18  3    New York      2010      19378102
19  4        Texas      2000      20851820
20  5        Texas      2010      25145561

```

在解决实际问题的時候，如果能將類似這樣的原始輸入數據的列直接轉換成 `MultiIndex`，通常將大有裨益。其實可以通過 `DataFrame` 的 `set_index` 方法實現，返回結果就會是一個帶多級索引的

`DataFrame`：

[117]:

```

1  pop_flat.set_index(['state', 'year'])
2  print(pop_flat)
3
4      state  year  population
5  0  California  2000      33871648
6  1  California  2010      37253956
7  2    New York  2000      18976457
8  3    New York  2010      19378102
9  4        Texas  2000      20851820
10 5        Texas  2010      25145561

```

在實踐中，我發現用這種重建索引的方法處理數據集非常好用。

3.6.5 多級索引的數據累計方法

前面已經介紹過一些 Pandas 自帶的數據累計方法，比如 `mean()`、`sum()` 和 `max()`。而對於層級索引數據，可以設置參數 `level` 實現對數據子集的累計操作。再一次以體檢數據為例：

[31]:

```

1  health_data

```

[31]:

	subject	Bob	Guido	Sue			
	type	HR	Temp	HR	Temp	HR	Temp
year	visit						
2013	1	47.0	35.7	24.0	37.3	35.0	35.9
	2	20.0	34.7	49.0	37.5	29.0	36.7
2014	1	41.0	38.3	34.0	36.0	37.0	37.3
	2	42.0	36.3	41.0	36.1	29.0	35.6

如果需要计算每一年各项指标的平均值，那么可以将参数 `level` 设置为索引 `year`：

[32]:

```
1 data_mean = health_data.mean(level='year')
2 data_mean
```

[32]:

subject	Bob	Guido	Sue			
type	HR	Temp	HR	Temp	HR	Temp
year						
2013	33.5	35.2	36.5	37.40	32.0	36.30
2014	41.5	37.3	37.5	36.05	33.0	36.45

如果再设置 `axis` 参数，就可以对列索引进行类似的累计操作了：

[33]:

```
1 data_mean.mean(axis=1, level='type')
```

[33]:

type	HR	Temp
year		
2013	34.000000	36.3
2014	37.333333	36.6

通过这两行数据，就可以获取每一年所有人的平均心率和体温了。这种语法其实就是 `GroupBy` 功能的快捷方式，将在 3.9 节详细介绍。尽管这只是一个简单的示例，但是其原理和实际工作中遇到的情况类似。

3.6.6 Panel数据

这里还有一些 Pandas 的基本数据结构没有介绍到，包括 `pd.Panel` 对象和 `pd.Panel4D` 对象。这两种数据结构可以分别看成是(一维数组) `Series` 和(二维数组) `DataFrame` 的三维与四维形式。如果熟悉 `Series` 和 `DataFrame` 的使用方法，那么 `Panel` 和 `Panel4D` 使用起来也会很简单，`ix`、`loc` 和 `iloc` 索引器(详情请参见 3.3 节)在高维数据结构上的用法更是完全相同。

但是本书并不打算进一步介绍这两种数据结构，我个人认为多级索引在大多数情况下都是更实用、更直观的高维数据形式。另外，`Panel` 采用密集数据存储形式，而多级索引采用稀疏数据存储形式。在解决许多真实的数据集时，随着维度的不断增加，密集数据存储形式的效率将越来越低。但是这类数据结构对一些有特殊需求的应用还是有用的。如果想对 `Panel` 与 `Panel4D` 数据结构有更多的认识，请参见 3.14 节。