



```
c    0.092908
d    0.281746
dtype: float64
```

层次化索引在数据重塑和基于分组的操作（如透视表生成）中扮演着重要的角色。例如，可以通过unstack方法将这段数据重新安排到一个DataFrame中：

```
In [16]: data.unstack()
Out[16]:
```

	1	2	3
a	-0.204708	0.478943	-0.519439
b	-0.555730	NaN	1.965781
c	1.393406	0.092908	NaN
d	NaN	0.281746	0.769023

unstack的逆运算是stack：

```
In [17]: data.unstack().stack()
Out[17]:
```

a	1	-0.204708
	2	0.478943
	3	-0.519439
b	1	-0.555730
	3	1.965781
c	1	1.393406
	2	0.092908
d	2	0.281746
	3	0.769023

dtype: float64

stack和unstack将在本章后面详细讲解。

对于一个DataFrame，每条轴都可以有分层索引：

```
In [18]: frame = pd.DataFrame(np.arange(12).reshape((4, 3)),
.....:                        index=[['a', 'a', 'b', 'b'], [1, 2, 1, 2]],
.....:                        columns=[['Ohio', 'Ohio', 'Colorado'],
.....:                                ['Green', 'Red', 'Green']])

In [19]: frame
Out[19]:
```

		Ohio		Colorado
		Green	Red	Green
a	1	0	1	2
	2	3	4	5
b	1	6	7	8
	2	9	10	11

各层都可以有名字（可以是字符串，也可以是别的Python对象）。如果指定了名称，它们就会显示在控制台输出中：

```
In [20]: frame.index.names = ['key1', 'key2']

In [21]: frame.columns.names = ['state', 'color']

In [22]: frame
Out[22]:
```

state		Ohio		Colorado
color		Green	Red	Green
key1	key2			
a	1	0	1	2
	2	3	4	5
b	1	6	7	8
	2	9	10	11

注意：小心区分索引名state、color与行标签。

```
In [23]: frame['Ohio']
Out[23]:
```

color	Green	Red
key1 key2		
a 1	0	1
2	3	4
b 1	6	7
2	9	10

```
MultiIndex.from_arrays([[ 'Ohio', 'Ohio', 'Colorado'], ['Green', 'Red', 'Green']],
                        names=[ 'state', 'color'])
```

有时，你需要重新调整某条轴上各级别的顺序，或根据指定级别上的值对数据进行排序。swaplevel接受两个级别编号或名称，并返回一个互换了级别的新对象（但数据不会发生变化）：

```
In [24]: frame.swaplevel('key1', 'key2')
Out[24]:
```

		Ohio	Colorado
state		Green	Red
color			Green
key2	key1		
1	a	0	2
2	a	3	5
1	b	6	8
2	b	9	11

```
In [25]: frame.sort_index(level=1)
Out[25]:
```

		Ohio		Colorado
color		Green	Red	Green
key1	key2			
a	1	0	1	2
b	1	6	7	8
a	2	3	4	5
b	2	9	10	11

```
In [26]: frame.swaplevel(0, 1).sort_index(level=0)
Out[26]:
```

		Ohio		Colorado
color		Green	Red	Green
key2	key1			
1	a	0	1	2
	b	6	7	8
2	a	3	4	5
	b	9	10	11

许多对DataFrame和Series的描述和汇总统计都有一个level选项，它用于指定在某条轴上求和的级别。再以上面那个DataFrame为例，我们可以根据行或列上的级别来进行求和：

```
In [27]: frame.sum(level='key2')
Out[27]:
state  Ohio      Colorado
```

```

color Green Red      Green
key2
1          6    8        10
2          12   14       16

In [28]: frame.sum(level='color', axis=1)
Out[28]:
color      Green    Red
key1 key2
a      1          2     1
      2          8     4
b      1         14     7
      2         20    10

```

这其实是利用了pandas的groupby功能，本书稍后将对其进行详细讲解。

## 使用DataFrame的列进行索引

人们经常想要将DataFrame的一个或多个列当做行索引来用，或者可能希望将行索引变成DataFrame的列。以下面这个DataFrame为例：

```

In [29]: frame = pd.DataFrame({'a': range(7), 'b': range(7, 0, -1),
.....:                        'c': ['one', 'one', 'one', 'two', 'two',
.....:                        'two', 'two'],
.....:                        'd': [0, 1, 2, 0, 1, 2, 3]})

In [30]: frame
Out[30]:
   a  b   c  d
0  0  7  one  0
1  1  6  one  1
2  2  5  one  2
3  3  4  two  0
4  4  3  two  1
5  5  2  two  2
6  6  1  two  3

```

DataFrame的set\_index函数会将其一个或多个列转换为行索引，并创建一个新的DataFrame：

```

In [31]: frame2 = frame.set_index(['c', 'd'])

In [32]: frame2
Out[32]:
      a  b
c  d
one 0  0  7
     1  1  6
     2  2  5
two 0  3  4
     1  4  3
     2  5  2
     3  6  1

```

默认情况下，那些列会从DataFrame中移除，但也可以将其保留下来：

```

In [33]: frame.set_index(['c', 'd'], drop=False)
Out[33]:
      a  b   c  d
c  d
one 0  0  7  one  0
     1  1  6  one  1
     2  2  5  one  2
two 0  3  4  two  0
     1  4  3  two  1
     2  5  2  two  2
     3  6  1  two  3

```

reset index的功能跟set index刚好相反，层次化索引的级别会被转移到列里面：

```
In [34]: frame2.reset_index()
Out[34]:
```

	c	d	a	b
0	one	0	0	7
1	one	1	1	6
2	one	2	2	5
3	two	0	3	4
4	two	1	4	3
5	two	2	5	2
6	two	3	6	1

## 8.2 合并数据集

pandas对象中的数据可以通过一些方式进行合并：

- `pandas.merge`可根据一个或多个键将不同DataFrame中的行连接起来。SQL或其他关系型数据库的用户对此应该会比较熟悉，因为它实现的就是数据库的join操作。
- `pandas.concat`可以沿着一条轴将多个对象堆叠到一起。
- 实例方法`combine_first`可以将重复数据拼接在一起，用一个对象中的值填充另一个对象中的缺失值。

我将分别对它们进行讲解，并给出一些例子。本书剩余部分的示例中将经常用到它们。

## ##数据库风格的DataFrame合并

数据集的合并（merge）或连接（join）运算都是通过一个或多个键将行连接起来的。这些运算是关系型数据库（基于SQL）的核心。pandas的merge函数是对数据应用这些算法的主要切入点。

以一个简单的例子开始：

```
In [35]: df1 = pd.DataFrame({'key': ['b', 'b', 'a', 'c', 'a', 'a', 'b'],
.....:                      'data1': range(7)})

In [36]: df2 = pd.DataFrame({'key': ['a', 'b', 'd'],
.....:                      'data2': range(3)})

In [37]: df1
Out[37]:
   data1 key
0      0   b
1      1   b
2      2   a
3      3   c
4      4   a
5      5   a
6      6   b

In [38]: df2
Out[38]:
   data2 key
0      0   a
1      1   b
2      2   d
```

这是一种多对一的合并。df1中的数据有多个被标记为a和b的行，而df2中key列的每个值则仅对应一行。对这些对象调用merge即可得到：

```
In [39]: pd.merge(df1, df2)
Out[39]:
```

	data1	key	data2
0	1	1	1
1	1	2	1
2	1	3	1
3	2	1	2
4	2	2	2
5	2	3	2
6	3	1	3
7	3	2	3
8	3	3	3

注意，我并没有指明要用哪个列进行连接。如果没有指定，merge就会将重叠列的列名当做键。不过，最好明确指定一下：

```
Out[40]:
```

	data1	key	data2
0	0	b	1
1	1	b	1
2	6	b	1
3	2	a	0
4	4	a	0
5	5	a	0

如果两个对象的列名不同，也可以分别进行指定：

```
Out[43]:
```

	data1	lkey	data2	rkey
0	0	b	1	b
1	1	b	1	b
2	6	b	1	b
3	2	a	0	a
4	4	a	0	a
5	5	a	0	a

可能你已经注意到了，结果里面c和d以及与之相关的数据消失了。默认情况下，merge做的是“内连接”；结果中的键是交集。其他方式还有“left”、“right”以及“outer”。外连接求取的是键的并集，组合了左连接和右连接的效果：

```
Out[44]:
```

	data1	key	data2
0	0.0	b	1.0
1	1.0	b	1.0
2	6.0	b	1.0
3	2.0	a	0.0
4	4.0	a	0.0
5	5.0	a	0.0
6	3.0	c	NaN
7	NaN	d	2.0

表8-1对这些选项进行了总结。

选项	说明
inner	使用两个表都有的键
left	使用左表中所有的键
right	使用右表中所有的键
outer	使用两个表中所有的键

表8-1 不同的连接类型

多对多的合并有些不直观。看下面的例子：

```
In [45]: df1 = pd.DataFrame({'key': ['b', 'b', 'a', 'c', 'a', 'b'],
.....:                      'data1': range(6)})

In [46]: df2 = pd.DataFrame({'key': ['a', 'b', 'a', 'b', 'd'],
.....:                      'data2': range(5)})

In [47]: df1
Out[47]:
  data1 key
0      0  b
1      1  b
2      2  a
3      3  c
4      4  a
5      5  b

In [48]: df2
Out[48]:
  data2 key
0      0  a
1      1  b
2      2  a
3      3  b
4      4  d

In [49]: pd.merge(df1, df2, on='key', how='left')
Out[49]:
  data1 key  data2
0      0  b      1.0
1      0  b      3.0
2      1  b      1.0
3      1  b      3.0
4      2  a      0.0
5      2  a      2.0
6      3  c      NaN
7      4  a      0.0
8      4  a      2.0
9      5  b      1.0
10     5  b      3.0
```

多对多连接产生的是行的笛卡尔积。由于左边的DataFrame有3个“b”行，右边的有2个，所以最终结果中就有6个“b”行。连接方式只影响出现在结果中的不同的键的值：

```
In [50]: pd.merge(df1, df2, how='inner')
Out[50]:
```

	data1	key	data2
0	0	b	1
1	0	b	3
2	1	b	1
3	1	b	3
4	5	b	1
5	5	b	3
6	2	a	0
7	2	a	2
8	4	a	0
9	4	a	2

要根据多个键进行合并，传入一个由列名组成的列表即可：

```
In [51]: left = pd.DataFrame({'key1': ['foo', 'foo', 'bar'],
.....:                      'key2': ['one', 'two', 'one'],
.....:                      'lval': [1, 2, 3]})

In [52]: right = pd.DataFrame({'key1': ['foo', 'foo', 'bar', 'bar'],
.....:                        'key2': ['one', 'one', 'one', 'two'],
.....:                        'rval': [4, 5, 6, 7]})

In [53]: pd.merge(left, right, on=['key1', 'key2'], how='outer')
Out[53]:
```

	key1	key2	lval	rval
0	foo	one	1.0	4.0
1	foo	one	1.0	5.0
2	foo	two	2.0	NaN
3	bar	one	3.0	6.0
4	bar	two	NaN	7.0

结果中会出现哪些键组合取决于所选的合并方式，你可以这样来理解：多个键形成一系列元组，并将其当做单个连接键（当然，实际上并不是这么回事）。

注意：在进行列－列连接时，DataFrame对象中的索引会被丢弃。

对于合并运算需要考虑的最后一个问题是对重复列名的处理。虽然你可以手工处理列名重叠的问题（查看前面介绍的重命名轴标签），但merge有一个更实用的suffixes选项，用于指定附加到左右两个DataFrame对象的重叠列名上的字符串：

```
In [54]: pd.merge(left, right, on='key1')
Out[54]:
```

	key1	key2_x	lval	key2_y	rval
0	foo	one	1	one	4
1	foo	one	1	one	5
2	foo	two	2	one	4
3	foo	two	2	one	5
4	bar	one	3	one	6
5	bar	one	3	two	7

```
In [55]: pd.merge(left, right, on='key1', suffixes=('_left', '_right'))
Out[55]:
```

	key1	key2_left	lval	key2_right	rval
0	foo	one	1	one	4
1	foo	one	1	one	5
2	foo	two	2	one	4
3	foo	two	2	one	5
4	bar	one	3	one	6
5	bar	one	3	two	7

merge的参数请参见表8-2。使用DataFrame的行索引合并是下一节的主题。



表8-2 merge函数的参数

参数	说明
left	参与合并的左侧DataFrame
right	参与合并的右侧DataFrame
how	“inner”、“outer”、“left”、“right”其中之一。默认为“inner”

参数	说明
on	用于连接的列名。必须存在于左右两个DataFrame对象中。如果未指定，且其他连接键也未指定，则以left和right列名的交集作为连接键
left_on	左侧DataFrame中用作连接键的列
right_on	右侧DataFrame中用作连接键的列
left_index	将左侧的行索引用作其连接键
right_index	类似于left_index
sort	根据连接键对合并后的数据进行排序，默认为True。有时在处理大数据集时，禁用该选项可获得更好的性能
suffixes	字符串值元组，用于追加到重叠列名的末尾，默认为('_', '_y')。例如，如果左右两个DataFrame对象都有“data”，则结果中就会出现“data_x”和“data_y”
copy	设置为False，可以在某些特殊情况下避免将数据复制到结果数据结构中。默认总是复制

indicator 添加特殊的列\_merge，它可以指明每个行的来源，它的值有left\_only、right\_only或both，根据每行的合并数据的来源。

## 索引上的合并

有时候，DataFrame中的连接键位于其索引中。在这种情况下，你可以传入left\_index=True或right\_index=True（或两个都传）以说明索引应该被用作连接键：

```
In [56]: left1 = pd.DataFrame({'key': ['a', 'b', 'a', 'a', 'b', 'c'],
.....:                        'value': range(6)})
```

```
In [57]: right1 = pd.DataFrame({'group_val': [3.5, 7]}, index=['a', 'b'])
```

```
In [58]: left1
Out[58]:
```

	key	value
0	a	0
1	b	1
2	a	2
3	a	3
4	b	4
5	c	5

```
In [59]: right1
Out[59]:
   group_val
a         3.5
b         7.0
```

```
In [60]: pd.merge(left1, right1, left_on='key', right_index=True)
```

```
Out[60]:
   key  value  group_val
0    a      0         3.5
2    a      2         3.5
3    a      3         3.5
1    b      1         7.0
4    b      4         7.0
```

由于默认的merge方法是求取连接键的交集，因此你可以通过外连接的方式得到它们的并集：

```
In [61]: pd.merge(left1, right1, left_on='key', right_index=True, how='outer')
```

```
Out[61]:
   key  value  group_val
0    a      0         3.5
2    a      2         3.5
3    a      3         3.5
1    b      1         7.0
4    b      4         7.0
5    c      5         NaN
```

对于层次化索引的数据，事情就有点复杂了，因为索引的合并默认是多键合并：

```
In [62]: lefth = pd.DataFrame({'key1': ['Ohio', 'Ohio', 'Ohio',
....:                                'Nevada', 'Nevada'],
....:                        'key2': [2000, 2001, 2002, 2001, 2002],
....:                        'data': np.arange(5.)})
```

```
In [63]: righth = pd.DataFrame(np.arange(12).reshape((6, 2)),
....:                          index=[['Nevada', 'Nevada', 'Ohio', 'Ohio',
....:                                'Ohio', 'Ohio'],
....:                          [2001, 2000, 2000, 2000, 2001, 2002]],
....:                          columns=['event1', 'event2'])
```

```
In [64]: lefth
Out[64]:
   data  key1  key2
0   0.0   Ohio  2000
1   1.0   Ohio  2001
2   2.0   Ohio  2002
3   3.0 Nevada  2001
4   4.0 Nevada  2002
```

```
In [65]: righth
Out[65]:
   event1  event2
Nevada  2001     0     1
        2000     2     3
Ohio    2000     4     5
        2000     6     7
```

这种情况下，你必须以列表的形式指明用作合并键的多个列（注意用`how='outer'`对重复索引值的处理）：

```
In [66]: pd.merge(leftth, rightth, left_on=['key1', 'key2'], right_index=True)
```

```
Out[66]:
```

	data	key1	key2	event1	event2
0	0.0	Ohio	2000	4	5
0	0.0	Ohio	2000	6	7
1	1.0	Ohio	2001	8	9
2	2.0	Ohio	2002	10	11
3	3.0	Nevada	2001	0	1

```
In [67]: pd.merge(lefth, righth, left_on=['key1', 'key2'],
.....:           right_index=True, how='outer')
```

```
Out[67]:
```

	data	key1	key2	event1	event2
0	0.0	Ohio	2000	4.0	5.0
0	0.0	Ohio	2000	6.0	7.0
1	1.0	Ohio	2001	8.0	9.0
2	2.0	Ohio	2002	10.0	11.0
3	3.0	Nevada	2001	0.0	1.0
4	4.0	Nevada	2002	NaN	NaN
4	NaN	Nevada	2000	2.0	3.0

同时使用合并双方的索引也没问题：

```
In [68]: left2 = pd.DataFrame([[1., 2.], [3., 4.], [5., 6.]],
.....:                        index=['a', 'c', 'e'],
.....:                        columns=['Ohio', 'Nevada'])
```

```
In [69]: right2 = pd.DataFrame([[7., 8.], [9., 10.], [11., 12.], [13, 14]],
.....:                        index=['b', 'c', 'd', 'e'],
.....:                        columns=['Missouri', 'Alabama'])
```

```
In [70]: left2
```

```
Out[70]:
```

	Ohio	Nevada
a	1.0	2.0
c	3.0	4.0
e	5.0	6.0

```
In [71]: right2
```

```
Out[71]:
```

	Missouri	Alabama
b	7.0	8.0
c	9.0	10.0
d	11.0	12.0
e	13.0	14.0

```
In [72]: pd.merge(left2, right2, how='outer', left_index=True, right_index=True)
```

```
Out[72]:
```

	Ohio	Nevada	Missouri	Alabama
a	1.0	2.0	NaN	NaN
b	NaN	NaN	7.0	8.0
c	3.0	4.0	9.0	10.0
d	NaN	NaN	11.0	12.0
e	5.0	6.0	13.0	14.0

DataFrame还有一个便捷的join实例方法，它能更为方便地实现按索引合并。它还可用于合并多个带有相同或相似索引的DataFrame对象，但要求没有重叠的列。在上面那个例子中，我们可以编写：

```
In [73]: left2.join(right2, how='outer')
```

```
Out[73]:
```

	Ohio	Nevada	Missouri	Alabama
a	1.0	2.0	NaN	NaN
b	NaN	NaN	7.0	8.0

c	3.0	4.0	9.0	10.0
d	NaN	NaN	11.0	12.0
e	5.0	6.0	13.0	14.0

因为一些历史版本的遗留原因，DataFrame的join方法默认使用的是左连接，保留左边表的行索引。它还支持在调用的DataFrame的列上，连接传递的DataFrame索引：

```
In [74]: left1.join(right1, on='key')
Out[74]:
```

	key	value	group_val
0	a	0	3.5
1	b	1	7.0
2	a	2	3.5
3	a	3	3.5
4	b	4	7.0
5	c	5	NaN

最后，对于简单的索引合并，你还可以向join传入一组DataFrame，下一节会介绍更为通用的concat函数，也能实现此功能：

```
In [75]: another = pd.DataFrame([[7., 8.], [9., 10.], [11., 12.], [16., 17.]],
.....:                          index=['a', 'c', 'e', 'f'],
.....:                          columns=['New York',
'Oregon'])
```

```
In [76]: another
Out[76]:
```

	New York	Oregon
a	7.0	8.0
c	9.0	10.0
e	11.0	12.0
f	16.0	17.0

```
In [77]: left2.join([right2, another])
Out[77]:
```

	Ohio	Nevada	Missouri	Alabama	New York	Oregon
a	1.0	2.0	NaN	NaN	7.0	8.0
c	3.0	4.0	9.0	10.0	9.0	10.0
e	5.0	6.0	13.0	14.0	11.0	12.0

```
In [78]: left2.join([right2, another], how='outer')
Out[78]:
```

	Ohio	Nevada	Missouri	Alabama	New York	Oregon
a	1.0	2.0	NaN	NaN	7.0	8.0
b	NaN	NaN	7.0	8.0	NaN	NaN
c	3.0	4.0	9.0	10.0	9.0	10.0
d	NaN	NaN	11.0	12.0	NaN	NaN
e	5.0	6.0	13.0	14.0	11.0	12.0
f	NaN	NaN	NaN	NaN	16.0	17.0

## 轴向连接

另一种数据合并运算也被称作连接（concatenation）、绑定（binding）或堆叠（stacking）。NumPy的concatenation函数可以用NumPy数组来做：

```
In [79]: arr = np.arange(12).reshape((3, 4))

In [80]: arr
Out[80]:
```

array([[	0,	1,	2,	3],	
	[	4,	5,	6,	7],
	[	8,	9,	10,	11]])

```
In [81]: np.concatenate([arr, arr], axis=1)
Out[81]:
```

array([[	0,	1,	2,	3,	0,	1,	2,	3],	
	[	4,	5,	6,	7,	4,	5,	6,	7],
	[	8,	9,	10,	11,	8,	9,	10,	11]])



在这个例子中，f和g标签消失了，是因为使用的是join='inner'选项。

你可以通过join\_axes指定要在其它轴上使用的索引：

```
In [91]: pd.concat([s1, s4], axis=1, join_axes=[['a', 'c', 'b', 'e']])
Out[91]:
   0    1
a  0.0  0.0
c  NaN  NaN
b  1.0  1.0
e  NaN  NaN
```

不过有个问题，参与连接的片段在结果中区分不开。假设你想要在连接轴上创建一个层次化索引。使用keys参数即可达到这个目的：

```
In [92]: result = pd.concat([s1, s1, s3], keys=['one', 'two', 'three'])

In [93]: result
Out[93]:
one    a    0
      b    1
two    a    0
      b    1
three  f    5
      g    6
dtype: int64

In [94]: result.unstack()
Out[94]:
      a    b    f    g
one   0.0  1.0  NaN  NaN
two   0.0  1.0  NaN  NaN
three NaN  NaN  5.0  6.0
```

如果沿着axis=1对Series进行合并，则keys就会成为DataFrame的列头：

```
In [95]: pd.concat([s1, s2, s3], axis=1, keys=['one', 'two', 'three'])
Out[95]:
      one  two  three
a   0.0  NaN   NaN
b   1.0  NaN   NaN
c   NaN  2.0   NaN
d   NaN  3.0   NaN
e   NaN  4.0   NaN
f   NaN  NaN   5.0
g   NaN  NaN   6.0
```

同样的逻辑也适用于DataFrame对象：

```
In [96]: df1 = pd.DataFrame(np.arange(6).reshape(3, 2), index=['a', 'b', 'c'],
.....:                      columns=['one', 'two'])

In [97]: df2 = pd.DataFrame(5 + np.arange(4).reshape(2, 2), index=['a', 'c'],
.....:                      columns=['three', 'four'])

In [98]: df1
Out[98]:
      one  two
a       0    1
b       2    3
c       4    5

In [99]: df2
Out[99]:
      three  four
a         5     6
c         7     8
```

```
In [100]: pd.concat([df1, df2], axis=1, keys=['level1', 'level2'])
Out[100]:
```

	level1		level2	
	one	two	three	four
a	0	1	5.0	6.0
b	2	3	NaN	NaN
c	4	5	7.0	8.0

如果传入的不是列表而是一个字典，则字典的键就会被当做keys选项的值：

```
In [101]: pd.concat({'level1': df1, 'level2': df2}, axis=1)
Out[101]:
```

	level1		level2	
	one	two	three	four
a	0	1	5.0	6.0
b	2	3	NaN	NaN
c	4	5	7.0	8.0

此外还有两个用于管理层次化索引创建方式的参数（参见表8-3）。举个例子，我们可以用names参数命名创建的轴级别：

```
In [102]: pd.concat([df1, df2], axis=1, keys=['level1', 'level2'],
Out[102]:
```

	level1		level2	
	one	two	three	four
a	0	1	5.0	6.0
b	2	3	NaN	NaN
c	4	5	7.0	8.0

最后一个关于DataFrame的问题是，DataFrame的行索引不包含任何相关数据：

```
In [103]: df1 = pd.DataFrame(np.random.randn(3, 4), columns=['a', 'b', 'c', 'd'])
In [104]: df2 = pd.DataFrame(np.random.randn(2, 3), columns=['b', 'd', 'a'])
In [105]: df1
Out[105]:
```

	a	b	c	d
0	1.246435	1.007189	-1.296221	0.274992
1	0.228913	1.352917	0.886429	-2.001637
2	-0.371843	1.669025	-0.438570	-0.539741

```
In [106]: df2
Out[106]:
```

	b	d	a
0	0.476985	3.248944	-1.021228
1	-0.577087	0.124121	0.302614

在这种情况下，传入ignore\_index=True即可：

```
In [107]: pd.concat([df1, df2], ignore_index=True)
Out[107]:
```

	a	b	c	d
0	1.246435	1.007189	-1.296221	0.274992
1	0.228913	1.352917	0.886429	-2.001637
2	-0.371843	1.669025	-0.438570	-0.539741
3	-1.021228	0.476985	NaN	3.248944
4	0.302614	-0.577087	NaN	0.124121









```
In [126]: result.unstack('state')
Out[126]:
```

state	Ohio	Colorado
number		
one	0	3
two	1	4
three	2	5

```
In [127]: s1 = pd.Series([0, 1, 2, 3], index=['a', 'b', 'c', 'd'])
In [128]: s2 = pd.Series([4, 5, 6], index=['c', 'd', 'e'])
In [129]: data2 = pd.concat([s1, s2], keys=['one', 'two'])
```

stack默认会滤除缺失数据，因此该运算是可逆的：

? æ°æ®è§æ 'i¼è□□å□□ã?□å□□å¹¶å□□é□□å;□.pdf[2020/7/14 18:20:06]

```
dtype: float64
```

在对DataFrame进行unstack操作时，作为旋转轴的级别将会成为结果中的最低级别：

```
In [135]: df = pd.DataFrame({'left': result, 'right': result + 5},
.....:                      columns=pd.Index(['left', 'right'], name='side'))
```

```
In [136]: df
Out[136]:
side      left  right
state  number
Ohio   one      0      5
       two      1      6
       three     2      7
Colorado one     3      8
        two     4      9
        three    5     10
```

```
In [137]: df.unstack('state')
Out[137]:
side left      right
state Ohio Colorado Ohio Colorado
number
one      0          3      5          8
two      1          4      6          9
three    2          5      7         10
```

当调用stack，我们可以指明轴的名字：

```
In [138]: df.unstack('state').stack('side')
Out[138]:
state      Colorado  Ohio
number side
one      left      3      0
        right     8      5
two      left      4      1
        right     9      6
three    left      5      2
        right    10      7
```

## 将“长格式”旋转为“宽格式”

多个时间序列数据通常是以所谓的“长格式”（long）或“堆叠格式”（stacked）存储在数据库和CSV中的。我们先加载一些示例数据，做一些时间序列规整和数据清洗：

```
In [139]: data = pd.read_csv('examples/macrodatab.csv')
```

```
In [140]: data.head()
Out[140]:
year quarter  realgdp  realcons  realinv  realgovt  realdpi  cpi  \
0  1959.0      1.0  2710.349   1707.4   286.898   470.045   1886.9  28.98
1  1959.0      2.0  2778.801   1733.7   310.859   481.301   1919.7  29.15
2  1959.0      3.0  2775.488   1751.8   289.226   491.260   1916.4  29.35
3  1959.0      4.0  2785.204   1753.7   299.356   484.052   1931.3  29.37
4  1960.0      1.0  2847.699   1770.5   331.722   462.199   1955.5  29.54
ml  tbilrate  unemp  pop  infl  realint
0  139.7      2.82   5.8  177.146  0.00   0.00
1  141.7      3.08   5.1  177.830  2.34   0.74
2  140.5      3.82   5.3  178.657  2.74   1.09
3  140.0      4.33   5.6  179.386  0.27   4.06
4  139.6      3.50   5.2  180.007  2.31   1.19
```

```
In [141]: periods = pd.PeriodIndex(year=data.year, quarter=data.quarter,
.....:                             name='date')
```

```
In [142]: columns = pd.Index(['realgdp', 'infl', 'unemp'], name='item')
```

```
In [143]: data = data.reindex(columns=columns)
In [144]: data.index = periods.to_timestamp('D', 'end')
In [145]: ldata = data.stack().reset_index().rename(columns={0: 'value'})
```

这就是多个时间序列（或者其它带有两个或多个键的可观察数据，这里，我们的键是date和item）的长格式。表中的每行代表一次观察。

关系型数据库（如MySQL）中的数据经常都是这样存储的，因为固定架构（即列名和数据类型）有一个好处：随着表中数据的添加，item列中的值的种类能够增加。在前面的例子中，date和item通常就是主键（用关系型数据库的说法），不仅提供了关系完整性，而且提供了更为简单的查询支持。有的情况下，使用这样的数据会很麻烦，你可能会更喜欢DataFrame，不同的item值分别形成一列，date列中的时间戳则用作索引。DataFrame的pivot方法完全可以实现这个转换：

```
In [147]: pivoted = ldata.pivot('date', 'item', 'value')
```

```
In [148]: pivoted
Out[148]:
item      infl      realgdp  unemp
date
1959-03-31  0.00      2710.349    5.8
1959-06-30  2.34      2778.801    5.1
1959-09-30  2.74      2775.488    5.3
1959-12-31  0.27      2785.204    5.6
1960-03-31  2.31      2847.699    5.2
1960-06-30  0.14      2834.390    5.2
1960-09-30  2.70      2839.022    5.6
1960-12-31  1.21      2802.616    6.3
1961-03-31 -0.40      2819.264    6.8
1961-06-30  1.47      2872.005    7.0
...
2007-06-30  2.75      13203.977    4.5
2007-09-30  3.45      13321.109    4.7
2007-12-31  6.38      13391.249    4.8
2008-03-31  2.82      13366.865    4.9
2008-06-30  8.53      13415.266    5.4
2008-09-30 -3.16      13324.600    6.0
2008-12-31 -8.79      13141.920    6.9
2009-03-31  0.94      12925.410    8.1
2009-06-30  3.37      12901.504    9.2
2009-09-30  3.56      12990.341    9.6
[203 rows x 3 columns]
```

前两个传递的值分别用作行和列索引，最后一个可选值则是用于填充DataFrame的数据列。假设有两个需要同时重塑的数据列：

```
In [149]: ldata['value2'] = np.random.randn(len(ldata))
```

```
In [150]: ldata[:10]
Out[150]:
   date      item      value      value2
0 1959-03-31  realgdp  2710.349  0.523772
1 1959-03-31    infl    0.000  0.000940
2 1959-03-31  unemp    5.800  1.343810
3 1959-06-30  realgdp  2778.801 -0.713544
4 1959-06-30    infl    2.340 -0.831154
5 1959-06-30  unemp    5.100 -2.370232
6 1959-09-30  realgdp  2775.488 -1.860761
7 1959-09-30    infl    2.740 -0.860757
8 1959-09-30  unemp    5.300  0.560145
9 1959-12-31  realgdp  2785.204 -1.265934
```

如果忽略最后一个参数，得到的DataFrame就会带有层次化的列：

```
In [151]: pivoted = ldata.pivot('date', 'item')
```

```
In [152]: pivoted[:5]
Out[152]:
item      value      value2
infl      realgdp unemp      infl      realgdp      unemp
```

```

date
1959-03-31    0.00   2710.349    5.8   0.000940   0.523772   1.343810
1959-06-30    2.34   2778.801    5.1  -0.831154  -0.713544  -2.370232
1959-09-30    2.74   2775.488    5.3  -0.860757  -1.860761   0.560145
1959-12-31    0.27   2785.204    5.6   0.119827  -1.265934  -1.063512
1960-03-31    2.31   2847.699    5.2  -2.359419   0.332883  -0.199543

```

```

In [153]: pivoted['value'][:5]
Out[153]:
item      infl    realgdp    unemp
date
1959-03-31    0.00   2710.349     5.8
1959-06-30    2.34   2778.801     5.1
1959-09-30    2.74   2775.488     5.3
1959-12-31    0.27   2785.204     5.6
1960-03-31    2.31   2847.699     5.2

```

注意，pivot其实就是用set\_index创建层次化索引，再用unstack重塑：

```
In [154]: unstacked = ldata.set_index(['date', 'item']).unstack('item')
```

```

In [155]: unstacked[:7]
Out[155]:
item      value      realgdp  unemp      value2      realgdp      unemp
date
1959-03-31    0.00   2710.349    5.8   0.000940   0.523772   1.343810
1959-06-30    2.34   2778.801    5.1  -0.831154  -0.713544  -2.370232
1959-09-30    2.74   2775.488    5.3  -0.860757  -1.860761   0.560145
1959-12-31    0.27   2785.204    5.6   0.119827  -1.265934  -1.063512
1960-03-31    2.31   2847.699    5.2  -2.359419   0.332883  -0.199543
1960-06-30    0.14   2834.390    5.2  -0.970736  -1.541996  -1.307030
1960-09-30    2.70   2839.022    5.6   0.377984   0.286350  -0.753887

```

## 将“宽格式”旋转为“长格式”

旋转DataFrame的逆运算是pandas.melt。它不是将一列转换到多个新的DataFrame，而是合并多个列成为一个，产生一个比输入长的DataFrame。看一个例子：

```

In [157]: df = pd.DataFrame({'key': ['foo', 'bar', 'baz'],
.....:                      'A': [1, 2, 3],
.....:                      'B': [4, 5, 6],
.....:                      'C': [7, 8, 9]})

```

```

In [158]: df
Out[158]:
   A  B  C  key
0  1  4  7  foo
1  2  5  8  bar
2  3  6  9  baz

```

key列可能是分组指标，其它的列是数据值。当使用pandas.melt，我们必须指明哪些列是分组指标。下面使用key作为唯一的分组指标：

```
In [159]: melted = pd.melt(df, ['key'])
```

```

In [160]: melted
Out[160]:
   key variable  value
0  foo         A      1
1  bar         A      2
2  baz         A      3
3  foo         B      4
4  bar         B      5
5  baz         B      6
6  foo         C      7
7  bar         C      8
8  baz         C      9

```

使用pivot, 可以重塑回原来的样子：

```
In [161]: reshaped = melted.pivot('key', 'variable', 'value')
In [162]: reshaped
Out[162]:
variable  A  B  C
key
bar      2  5  8
baz      3  6  9
foo      1  4  7
```

因为pivot的结果从列创建了一个索引, 用作行标签, 我们可以使用reset\_index将数据移回列：

```
In [163]: reshaped.reset_index()
Out[163]:
variable key  A  B  C
0      bar  2  5  8
1      baz  3  6  9
2      foo  1  4  7
```

你还可以指定列的子集, 作为值的列：

```
In [164]: pd.melt(df, id_vars=['key'], value_vars=['A', 'B'])
Out[164]:
   key variable  value
0  foo         A       1
1  bar         A       2
2  baz         A       3
3  foo         B       4
4  bar         B       5
5  baz         B       6
```

pandas.melt也可以不用分组指标：

```
In [165]: pd.melt(df, value_vars=['A', 'B', 'C'])
Out[165]:
   variable  value
0         A       1
1         A       2
2         A       3
3         B       4
4         B       5
5         B       6
6         C       7
7         C       8
8         C       9
```

```
In [166]: pd.melt(df, value_vars=['key', 'A', 'B'])
Out[166]:
   variable  value
0        key    foo
1        key    bar
2        key    baz
3         A       1
4         A       2
5         A       3
6         B       4
7         B       5
8         B       6
```

## #8.4 总结

现在你已经掌握了pandas数据导入、清洗、重塑, 我们可以进一步学习matplotlib数据可视化。我们在稍后会回到pandas, 学习更高级的分析。