

Python 数据分析与数据挖掘（ Python for Data Analysis&Data Mining ）

Chap 4 pandas数据探索 Data Exploration

内容：

- 数据质量分析：缺失、异常、不一致、重复或特殊符号
- 数据特征分析：分布分析、汇总和描述统计分析、相关性分析

实践：

- pandas的Sries和DataFrame
- pandas进行数据探索

这节课是学习pandas对数据的各种分析，目的是了解数据。下节课对数据的内容进行预处理。后面的课程对数据进行挖掘，从目前已有数据中挖掘出知识，用于预测新的数据。

pandas介绍

pandas是基于NumPy构建的，含有使数据分析工作变得更快更简单的高级数据结构和操作工具。从2008年后，pandas逐渐成长为一个非常大的库，能解决越来越多的数据分析问题。

pandas名字源于panel data（面板数据，是计量经济学中关于多维结构化数据集的一个术语）以及Python data analysis（Python数据分析），很适合金融数据分析应用的工具

- 具有高性能的数组计算功能以及电子表格和关系型数据库（如SQL）灵活的数据处理能力
- 提供了复杂精细的索引功能，更便捷地完成重塑、切片和切块、聚合以及选取数据子集等操作
- 对于金融行业的用户，pandas提供了大量适用于金融数据的高性能时间序列功能和工具

我们使用下面的pandas引入约定：

```
from pandas import Series, DataFrame  
import pandas as pd
```

因为Series和DataFrame用的次数非常多，因此将其引入本地命名空间中会更方便

准备工作：导入库，配置环境等

```
In [ ]: from __future__ import division
import os, sys

# 导入库并为库起个别名
import numpy as np
import pandas as pd
from pandas import Series, DataFrame

# 启动绘图
%matplotlib inline
import matplotlib.pyplot as plt

# 常用全局配置
np.random.seed(12345)
np.set_printoptions(precision=4)
plt.rc('figure', figsize=(10, 6))
br = '\n'
```

pandas的数据结构

pandas中两个最主要的数据结构就是：Series 和 DataFrame，为大多数应用提供了一种可靠的、易于使用的基础。

1. Series

Series是一种类似于一维数组的对象，它由一组数据（各种NumPy数据类型）以及一组与之相关的数据标签（即索引）组成。如下，仅由一组数据即可产生最简单的Series：

- Series的字符串表现形式为：索引在左边，值在右边
- 如果没有为数据指定索引，会自动创建一个0到N-1的（N为数据的长度）的整数型索引。

```
In [ ]: # 由列表创建一个Series序列
s = Series([4, 7, -5, 3]) # 没有指定索引
#s = Series([4, 7, -5, 3], index=['1', '2', '3', '4']) # 指定索引
s
```

```
In [ ]: # 可以通过Series的values和index属性获取其数组表示形式和索引对象:
print s.values, br
print s.index, br
```

```
In [ ]: # 希望所创建的Series带有一个可以对各个数据点进行标记的索引。
# 所以，Series可以看做是一个定长的有序字典。
s2 = Series([4, 7, -5, 3], index=['d', 'b', 'a', 'c']) # 指定索引
print s2.values, br
print s2.index, br
```

```
In [ ]: # 与普通NumPy数组相比，可以通过索引的方式选取Series中的单个或一组值:
s2['a']
```

```
In [ ]: # 修改d索引的值
print s2['d'], br
s2['d'] = 6
s2
```

```
In [ ]: # 按索引切片数据  
print s2[['c', 'a', 'd']]
```

- NumPy数组运算（如根据布尔型数组进行过滤、标量乘法、应用数学函数等）都会保留索引和值之间的链接：

```
In [ ]: # 数据过滤、标量乘法等数组运算都会保留索引和值之间的链接:  
s2[s2 > 0]
```

```
In [ ]: s2 * 2
```

```
In [ ]: np.exp(s2)
```

- 还可以把Series看成是一个定长的有序字典，因为它是索引值到数据值的一个映射。它可以用在许多原本需要字典参数的函数中：

```
In [ ]: print 'b' in s2  
print 'e' in s2
```

```
In [ ]: # 可以直接通过字典来创建Series: # sdata是Python的字典对象  
sdata = {'Ohio': 35000, 'Texas': 71000, 'Oregon': 16000, 'Utah': 5000}  
s3 = Series(sdata) # Series中的索引就是原字典的键key (有序排列)  
  
print s3, br  
print s3.index, s3.values
```

在下面这个例子中，sdata跟states索引相匹配的那3个值会被找出来并放到相应的位置上，但由于‘California’所对应的sdata值找不到，所以其结果就是NaN（即“非数字”，not a number，在pandas中它用于表示缺失或NA值）。

```
In [ ]: states = ['California', 'Ohio', 'Oregon', 'Texas']  
print sdata, br  
# 寻找跟states索引相匹配的那3个值会被找出来并放到相应的位置上，创建Series序列  
s4 = Series(sdata, index=states)  
s4
```

```
In [ ]: # pandas 的isnull用于检测缺失数据missing data, 返回布尔型数组的值  
pd.isnull(s4)
```

```
In [ ]: # pandas的notnull用于检测非空数据  
pd.notnull(s4)
```

Series也有类似的实例方法：

- obj.isnull() 等同于 pd.isnull(obj)
- obj.notnull() 等同于 pd.notnull(obj)

```
In [ ]: # Series的isnull()方法也可以用于检测缺失数据missing data, 返回布尔型数组的值  
s4.isnull()
```

```
In [ ]: # Series的notnull()方法也可以用于检测非空数据  
s4.notnull()
```

```
In [ ]: # Series 最重要的一个功能是，它在算术运算中会自动对齐不同索引的数据  
print s3, '\n', s4  
# Series在算术运算中自动对齐不同索引的数据  
s3 + s4
```

```
In [ ]: # Series对象本身及其索引都有一个name属性，该属性跟pandas其他的关键功能关系非常密切  
s4.name = 'population'  
s4.index.name = 'state'  
s4
```

```
In [ ]: # Series的索引可以通过赋值的方式就地修改：  
print s, br  
print s.index, br  
s.index = ['Bob', 'Steve', 'Jeff', 'Ryan']  
s
```

2. DataFrame

DataFrame是一个表格型的数据结构，它含有一组有序的列，每列可以是不同的值类型（数值、字符串、布尔值等）。

DataFrame既有行索引也有列索引，它可以被看做由Series组成的字典（共用同一个索引，如序列号，可以认为是行索引）。每列可以看作一个Series。

跟其他类似的数据结构相比（如R的data.frame），DataFrame中面向行和面向列的操作基本上是平衡的。其实，DataFrame中的数据是以一个或多个二维块存放的（而不是列表、字典或别的一维数据结构）。

注意：虽然DataFrame是以二维结构保存数据的，仍然可以轻松地将其表示为更高维度的数据（层次化索引的表格型结构，这是pandas中许多高级数据处理功能的关键要素）。

构建DataFrame的方法有很多：（1）最常用的一种是直接传入一个由等长列表或NumPy数组组成的字典。结果DataFrame会自动加上索引（跟Series一样），且全部列会被有序排列。

```
In [ ]: # data是字典，每个key值对应的value是一个list，每个key值和value list对应一列  
data = {'state': ['Ohio', 'Ohio', 'Ohio', 'Nevada', 'Nevada'],  
        'year': [2000, 2001, 2002, 2001, 2002],  
        'pop': [1.5, 1.7, 3.6, 2.4, 2.9]}  
frame = DataFrame(data) # DataFrame会自动加上索引（跟Series一样），且全部列会被有序排列  
frame
```

```
In [ ]: # 如果指定了列序列，则DataFrame的列就会按照指定顺序进行排列。  
DataFrame(data, columns=['year', 'state', 'pop'])
```

```
In [ ]: # 与Series一样，如果传入的列在数据中找不到，就会产生NA值：  
frame2 = DataFrame(data, columns=['year', 'state', 'pop', 'debt'],  
                    index=['one', 'two', 'three', 'four', 'five'])  
frame2 # columns是列索引，index是行索引
```

通过类似字典标记的方法或属性的方式，可以将DataFrame的列获取为一个Series，以下两种方式结果相同，都是将DataFrame的列取出为一个Series：

- 1) 字典标记的方式 frame['state']
- 2) 属性的方式 frame.state

注意：返回的Series用于原DataFrame相同的索引，且其name属性也已经被相应地设置好了。

每个列可以看作是一个Series，即column:Series，其中column是列索引

```
In [ ]: print frame2.columns # 返回所有的列索引  
        print frame2.index # 返回所有的行索引
```

```
In [ ]: frame2['state'] # 使用列索引返回指定列，对行操作无效  
# frame2.state # 相同
```

```
In [ ]: frame2.year # 返回指定列
```

```
In [ ]: # 返回行只能通过位置或名称的方式进行获取，比如用索引字段ix获取某个行：  
print frame2.index # 返回所有的行索引  
frame2.ix['three'] # 返回某特定索引的行数据
```

```
In [ ]: # 列可以通过赋值的方式进行就地修改。下面的方式对整个列的所有行的值都进行了就地修改  
print frame2, br  
frame2['debt'] = 16.5 # 给空的“debt”列附上一个标量值或一组值，就地修改  
frame2
```

```
In [ ]: frame2['debt'] = np.arange(5.) # 给debt列附上一组值，就地修改  
frame2
```

注意：将列表或数组赋值给某个列时，其长度必须跟DataFrame的长度相匹配。如果赋值的是一个Series，就会精确匹配DataFrame的索引，所有的空位都将被填上缺失值：

```
In [ ]: val = Series([-1.2, -1.5, -1.7], index=['two', 'four', 'five'])  
print frame2  
frame2['debt'] = val # 精确匹配DataFrame的索引，所有的空位都将被填上缺失值  
frame2 # 原来有值的位置也被填上了缺失值，说明Series整体覆盖原值
```

```
In [ ]: # 获取DataFrame表格中的值  
print frame2.ix['one']['state'] # 先按行索引，再按列索引  
print frame2['debt'][1] # 先按列索引，再按行索引
```

```
In [ ]: # 按布尔值增加新的列  
frame2['eastern'] = frame2.state == 'Ohio'  
frame2
```

```
In [ ]: # 关键词del用于删除列  
del frame2['eastern']  
frame2.columns
```

注意：通过索引方式返回的列是相应数据的视图，并不是副本。因此，对返回的Series所做的任何就地修改全部都会反映到源DataFrame上。通过Series的copy方法即可显式地复制列。

另一种常见的数据形式是嵌套字典（也就是字典的字典）。如果把它传给DataFrame，它就会被解释为：外层字典的键作为列，内层键则作为行索引。

```
In [ ]: # 嵌套字典（也就是字典的字典）
pop = {'Nevada': {2001: 2.4, 2002: 2.9},
       'Ohio': {2000: 1.5, 2001: 1.7, 2002: 3.6}}
pop
```

```
In [ ]: # 把嵌套字典传给DataFrame: 外层字典的键作为列，内层键则作为行索引。
frame3 = DataFrame(pop)
frame3
```

```
In [ ]: # 可以对DataFrame进行转置
frame3.T
```

```
In [ ]: # 内层字典的键会被合并、排序以形成最终的索引。如果显式指定了索引，则不会合并。
print pop
DataFrame(pop, index=[2001, 2002, 2003])
```

```
In [ ]: # 由Series组成的字典差不多也是一样的用法
print frame3, br
print frame3['Ohio'][:-1]
pdata = {'Ohio': frame3['Ohio'][:-1],
          'Nevada': frame3['Nevada'][:2]}
DataFrame(pdata)
```

```
In [ ]: # 设置frame的index和column的name属性，则会显式出来行和列的名字
print frame3.index.name
frame3.index.name = 'year'; frame3.columns.name = 'state'
frame3
```

```
In [ ]: # 跟Series一样，values属性也会以二维ndarray的形式返回DataFrame中的数据
frame3.values
```

```
In [ ]: # 如果DataFrame的各列的数据类型不同，则值数组的数据类型就会选用能兼容所有列的数据类型。
frame2.values
```

下面的表列出了DataFrame构造函数所能接受的各种数据

表：可以输入给DataFrame构造器的数据

类型	说明
二维ndarray	数据矩阵，还可以传入行标和列标
由数组、列表或元祖组成的字典	每个序列会变成DataFrame中的一列，所有序列的长度必须相同
NumPy的结构化/记录数组	类似于“由数组组成的字典”
由Series组成的字典	每个Series会成为一列。如果没有显式索引，则各Series的索引会被合并成结果的行索引
由字典组成的字典	各内层字典会成为一列。键会被合并成结果的行索引，跟“由Series组成的字典”的情况一样
字典或Series的列表	各项将会成为DataFrame的一行。字典键或Series索引的并集将会成为DataFrame的列标
由列表或元组组成的列表	类似于“二维ndarray”
另一个DataFrame	该DataFrame的索引将会被沿用，除非显式指定了其他索引
NumPy的MaskedArray	类似于“二维ndarray”的情况，只是掩码值在结果DataFrame会变成NA或缺失值。

索引对象 (Index objects)

pandas的索引对象负责管理轴标签和其他元数据（比如轴名称等）。构建Series或DataFrame时，所用到的任何数组或其他序列的标签都会被转换成一个Index：

- Index对象是不可修改的（immutable），因此用户不能对其进行修改。这个属性非常重要，因为这样才可以使得Index对象在多个数据结构之间实现安全共享。
- Index不可修改的属性非常重要，因为这样才能使得Index对象在多个数据结构之间安全共享

```
In [ ]: obj = Series(range(3), index=['a', 'b', 'c'])
index = obj.index
index
```

```
In [ ]: index[1:]
```

```
In [ ]: # index是不能修改的，下面的语句会报错，Index对象是不可修改的 (immutable)
#index[1] = 'd'
```

```
In [ ]: index = pd.Index(np.arange(3))
s2 = Series([1.5, -2.5, 0], index=index)
s2.index is index
```

pandas中主要的Index对象

下面的表列出了pandas库中内置的Index类。Index甚至可以被继承从而实现特别的轴索引功能。

类	说明
Index	最泛化的Index对象，将轴标签表示为一个由Python对象组成的NumPy数组
Int64Index	针对整数的特殊Index
MultIndex	“层次化”索引对象，表示单个轴上的多层次索引。可以看做由元组组成的数组
DatetimeIndex	存储纳秒级时间戳（用NumPy的datetime64类型表示）
PeriodIndex	针对Period数据（时间间隔）的特殊Index

除了长得像数组，Index的功能也类似一个固定大小的集合。

- 每个索引都有一些方法和属性，可用于设置逻辑并回答有关该索引所包含数据的常见问题。
- 索引类似集合，集合的方法和属性也基本可以适用Index对象。

Index的方法和属性

方法	说明
append	连接另一个Index对象，产生一个新的Index
diff	计算差集，并得到一个Index
intersection	计算交集
union	计算并集
in/is	计算一个索引各值是否都包含着参数集合中的布尔型数组
delete	删除索引 i 处的元素，并得到新的Index
drop	删除传入的值，并得到新的Index
insert	将元素插入到索引 i 处，并得到新的Index
is_monotonic	当各元素均大于等于前一个元素时，返回True
is_unique	当Index没有重复值时，返回True
unique	计算Index中唯一值的数组

```
In [ ]: # 计算一个索引各值是否包含参数集合中数值的布尔型数组
print frame3
print 'Ohio' in frame3.columns # 是否在数据的列索引名中
print 2.4 in frame3.values # 是否在数据的值中
print 2003 in frame3.index # 是否在数据的行索引名中
```

pandas的基本功能 (Essential functionality)

下面介绍操作Series和DataFrame中的数据的基本手段，来探索pandas在数据分析和处理方面的功能。

重新索引 (Reindexing)

pandas对象的一个重要方法是 reindex，其作用是创建一个适应新索引的新对象。

- 调用该Series的reindex会根据新索引进行重排。如果某个索引值不存在，就引入缺失值
- fill_value=0 对missing value进行填充，缺失值填充值=0

```
In [ ]: s = Series([4.5, 7.2, -5.3, 3.6], index=['d', 'b', 'a', 'c']) # 创建序列s  
s2 = s.reindex(['a', 'b', 'c', 'd', 'e']) # 对序列s重排索引得到s2, 默认缺失值填充NaN  
s2
```

```
In [ ]: # fill_value=0 对missing value进行填充，缺失值填充值=0  
s.reindex(['a', 'b', 'c', 'd', 'e'], fill_value=0)
```

- 对于时间序列这样的有序数据，重新索引时需要做插值处理，method选项可达到这个目的。
- reindex()的参数 method='ffill' 表示可以实现前向值填充，将前一个值传播给下一个元素，propagate last valid observation forward to next valid

```
In [ ]: s3 = Series(['blue', 'purple', 'yellow'], index=[0, 2, 4])  
s3.reindex(range(6), method='ffill') # 前向值填充，即将前一个值传播给下一个元素
```

下面的表格列出了reindex可用的method选项。

有时我们可能需要比前向和后向填充更为精准的插值方式。

```
Series.reindex(range(6), method='ffill')
```

```
method= {None, 'backfill'/'bfill', 'pad'/'ffill', 'nearest'}
```

参数	说明
default	不填充 (don't fill gaps)
ffill或pad	前向填充 (或搬运) 值(propagate last valid observation forward to next valid)
bfill或backfill	后向填充 (或搬运) 值 (use next valid observation to fill gap)
nearest	使用最近的值填充 (use nearest valid observations to fill gap)

- 对于DataFrame，reindex可以修改（行）索引、列，或两个都修改。
- 如果仅传入一个序列，则会重新索引行。
- reindex并没有就地修改，即，没有改变源数据

```
In [ ]: frame = DataFrame(np.arange(9).reshape((3, 3)), index=['a', 'c', 'd'],  
columns=['Ohio', 'Texas', 'California'])  
frame
```

```
In [ ]: # 只传入一个索引序列，则reindex只修改行索引，即添加一个新行  
frame2 = frame.reindex(['a', 'b', 'c', 'd']) # 默认添入新行的值为NaN  
frame2
```

```
In [ ]: # 使用columns关键字，reindex可以重新索引列，即实现修改列  
states = ['Texas', 'Utah', 'California']  
frame.reindex(columns=states)
```

```
In [ ]: # 也可以同时对行和列进行重新索引，而插值则只能按行应用（即轴0）  
frame.reindex(index=['a', 'b', 'c', 'd'], method='ffill',  
columns=states)
```

```
In [ ]: # 注意：reindex并没有修改源数据  
frame
```

```
In [ ]: # 利用ix的标签索引功能，重索引任务可以变得更简洁  
frame.ix[['a', 'b', 'c', 'd'], states]
```

下表列出reindex函数的参数和说明

参数	说明
index	用作索引的新序列。即可以是Index实例，也可以是其他序列型的Python数据结构。Index会被完全使用，就像没有任何复制一样
method	插值（填充）方式，具体可以参见前面的表格
fill_value	在重新索引的过程中，需要引入缺失值时使用的替代值
limit	前向或后向填充时的最大填充量
level	在MultiIndex的指定级别上匹配简单索引，否则选取其子集
copy	默认为True，无论如何都复制；如果为False，则新旧相等就不复制

删除指定轴上的项 (Dropping entries from an axis)

删除某条轴上的一个或多个项很简单，只要一个索引数组或列表即可。由于需要执行一些数据整理和集合逻辑，所以drop方法返回的是一个在指定轴上删除了指定值的新对象(没有改变源数据，即，没有对源数据删除)。

```
In [ ]: # 对series序列，只能删除指定行上的值  
s = Series(np.arange(5.), index=['a', 'b', 'c', 'd', 'e'])  
new_s = s.drop('c') # 丢弃指定的c行  
new_s
```

```
In [ ]: s.drop(['d', 'c']) # 删除指定的d和c行
```

```
In [ ]: # 对于DataFrame，可以删除任意轴上的索引值。axis = 1即删除列轴  
data = DataFrame(np.arange(16).reshape((4, 4)),  
                 index=['Ohio', 'Colorado', 'Utah', 'New York'],  
                 columns=['one', 'two', 'three', 'four'])  
data
```

```
In [ ]: data.drop('Colorado') # 删除指定的索引值，axis默认为0，即行轴
```

```
In [ ]: data.drop(['Colorado', 'Ohio']) # 删除指定的两个索引值, axis默认为0, 即行轴
```

```
In [ ]: data.drop('two', axis=1) # 删除指定的列索引值, 必须指定axis=1, 即列轴
```

```
In [ ]: data.drop(['two', 'four'], axis=1)
```

```
In [ ]: #data.drop('two') # 报错: 删除指定的列索引值, 必须增加axis=1, 即列轴
```

索引、选取和过滤 (Indexing, selection, and filtering)

Series索引 (obj[...]) 的工作方式类似于NumPy数组的索引，只不过Series的索引值不只是整数。

```
In [ ]: s = Series(np.arange(4.), index=['a', 'b', 'c', 'd'])
print s['b'] # 直接返回b索引对应的value
print s[['b']] # 返回b和c索引对应的行, 即 key 和 item的value
```

```
In [ ]: print s[1] # series的索引只有行索引, 可以使用整数索引代替索引值
print s[2:4] # 选取片段
print s[['b', 'a', 'd']] # 对多个行进行索引, 需要使用list[]指明索引
print s[[1, 3]] # 选取多个行
```

```
In [ ]: s[s < 2] # 按布尔值选取
```

- 利用标签的切片运算与普通的Python切片运算不同，其末端是包含的闭区间 (inclusive)

```
In [ ]: s['b':'c'] #利用标签的切片与普通的Python不同, 末端是包含闭区间, 比较s[1:2]
```

```
In [ ]: s['b':'c'] = 5 # 给b和c索引的项赋值, 改变源数据的值
s
```

DataFrame中获取数据，无论是行切片或者列切片，都需要进行数据切片。

要从DataFrame选取列：

- 直接使用标签索引，获取一个或多个列，`data['column1']`或`data[['column2','column3']]`或`data.column1`

要从DataFrame选取行的方式有三种：

- 1) 使用数字索引，获取一个或多个行，`data[2:4]`(选取第3和第4行数据),即使一行，也需要使用切片`[i:i+1]`获取第i行
- 2) 通过布尔型数组选取行 `data[data[column] > 5]`
- 3) 为了在行上进行标签索引，引入专门的索引字段 `ix`，可以通过NumPy式的标记法以及轴标签从 DataFrame中选取行和列的子集。这也是一种重新索引的简单手段。

索引字段`ix`可以按照标签索引选取行和列的子集 `data. ix[' Colorado', [' two', ' three']]`

根据整数位置选取单列或单行，并返回一个Series：

- `print data.iloc[:,1]` 根据整数位置选取单列，并返回一个Series
- `print data.iloc[3]` 根据整数位置选取单行，并返回一个Series

根据标签选取单行或单列，并返回一个Series

- `print data.xs('Utah')`
- `print data.xs('one', axis=1)`

```
In [ ]: data = DataFrame(np.arange(16).reshape((4, 4)),
                      index=['Ohio', 'Colorado', 'Utah', 'New York'],
                      columns=['one', 'two', 'three', 'four'])
data
```

```
In [ ]: data['two'] # 列索引就是获取一个列，进行切片
data.two # 同上，获取一个列，进行切片
data[['three', 'one']] # 获取多个列，对列索引使用list[]进行切块
```

```
In [ ]: # 要从DataFrame选取行，使用切片方式
print data[0:1], br # 使用数字索引，默认就是行索引（轴为0）的获取行
print data[:2], br #
print data[1:2] # 如果只获取一行，不能直接使用这行的数字索引，而要使用[i:i+1]获取第i行数据
```

```
In [ ]: #print data, br
data[data['three'] > 5] # 对DataFrame也可以通过布尔型数组选取数据
```

```
In [ ]: print data, br
print data < 5 # 由标量比较运算获得布尔型DataFrame
```

```
In [ ]: data[data < 5] = 0 # 利用标量比较运算得到的布尔型DataFrame进行索引
data
```

```
In [ ]: data.ix['Colorado', ['two', 'three']] # 索引字段ix可以按照标签索引选取行和列的子集
data.ix[['Colorado', 'Utah'], [3, 0, 1]] # 索引字段ix可以按照标签索引选取行和列的子集
data.ix[2] # 返回行索引（轴=0），第3行
data.ix[2][1] # 返回第3行的第2个项，采用切块的方法取单个元素项
```

```
In [ ]: data.ix['Colorado', ['two', 'three']] # 索引字段ix可以按照标签索引选取行和列的子集
          data.ix[['Colorado', 'Utah'], [3, 0, 1]] # 索引字段ix可以按照标签索引选取行和列的子集
          data.ix[2] # 返回行索引(轴=0), 第3行
          data.ix[2][1] # 返回第3行的第2个项, 采用切块的方法取单个元素项
          data.ix[:'Utah', 'two'] # 切块操作取得单个元素, 包含Utah的闭区间之前的行, 列=two
          data.ix[data.three > 5, :3] # 结合比较运算切块
```

```
In [ ]: print data, br
          # 根据整数位置选取单列或单行, 并返回一个Series
          data.iloc[:, 1], br # 根据整数位置选取单列, 并返回一个Series
          data.iloc[3], br # 根据整数位置选取单行, 并返回一个Series
          # 根据标签选取单行或单列, 并返回一个Series
          data.xs('Utah'), br # 根据标签选取单列, 默认 axis=0
          data.xs('one', axis=1), br # 根据标签选取单列, axis=1
```

下面的表格列出了针对DataFrame数据的选取和重排方式

对pandas对象中的数据的选取和重排方式有很多，在使用层次化索引时还能用到一些别的方法。

对于pandas对象，必须输入obj[:, col]才能选取列实在有些啰嗦，而且很容易出错，因为列的选取是一种最常见的操作，因此，在设计pandas的时候就把所有的标签索引功能都放到 ix 中了。如下表显示的选项。

下表列出DataFrame的索引选项

类型	说明
obj.val	获取DataFrame的单个列，返回的一个Series
obj[val]	选取DataFrame的单个列或一组列。在一些特殊情况下会比较便利：布尔型数组（过滤行）、切片（行切片）、布尔型DataFrame（根据条件设置值）
obj.ix[val]	选取DataFrame的单个行或一组行
obj.ix[:, val]	选取单个列或列子集
obj.ix[val1, val2]	同时选取行和列（切块操作）
reindex方法	将一个或多个轴匹配到新索引
xs 方法	根据标签选取单行或单列，并返回一个Series
icol、irow方法	根据整数位置选取单列或单行，并返回一个Series
get_value、set_value方法	根据行标签和列标签选取/设置单个值

算术运算和数据对齐 (Arithmetic and data alignment)

- pandas最重要的一个功能是，它可以对不同索引的对象进行算术运算。在将对象相加时，如果存在不同的索引对，则结果的索引就是该索引对的并集。
- 自动的数据对齐操作中不重叠的索引处引入了NA值。缺失值会在算术运算过程中传播。
- 对于DataFrame，对齐操作会同时发生在行和列上。

```
In [ ]: s1 = Series([7.3, -2.5, 3.4, 1.5], index=['a', 'c', 'd', 'e'])
          s2 = Series([-2.1, 3.6, -1.5, 4, 3.1], index=['a', 'c', 'e', 'f', 'g'])
          s1 + s2 # 自动数据对齐, 在不重叠的索引处引入NA值
```

```
In [ ]: list('bcd') # 把字符串转为list
df1 = DataFrame(np.arange(9.).reshape((3, 3)), columns=list('bcd'),
               index=['Ohio', 'Texas', 'Colorado'])
df2 = DataFrame(np.arange(12.).reshape((4, 3)), columns=list('bde'),
               index=['Utah', 'Ohio', 'Texas', 'Oregon'])
# 对于DataFrame对齐操作会同时发生在行和列上, 相加后返回一个新的DataFrame,
# 其索引和列为原来那两个DataFrame的并集。
df1 + df2
```

在算术方法中填充值 (Arithmetic methods with fill values)

- 在对不同索引的对象进行算术运算时，可能希望当一个对象中某个轴标签中另一个对象中找不到时填充一个特殊值（比如0），而不是默认的NA值。
- 使用df1的add方法，传入df2以及一个fill_value参数, fill_value=0 表示对于找不到的对象时默认值为0，
df1.add(df2, fill_value=0)

```
In [ ]: df1 = DataFrame(np.arange(12.).reshape((3, 4)), columns=list('abcd'))
df2 = DataFrame(np.arange(20.).reshape((4, 5)), columns=list('abcde'))
df1 + df2 # 将它们相加, 没有重叠的位置就会产生NA值
df1.add(df2, fill_value=0) # 使用df1的add方法, 传入df2以及一个fill_value参数
```

```
In [ ]: # 类似的, 在对Series或DataFrame重新索引时, 也可以指定一个填充值:
df1.reindex(columns=df2.columns, fill_value=0)
```

下表列出灵活的算术方法

方法	说明
add	用于加法 (+) 的方法
sub	用于减法 (-) 的方法
div	用于除法 (/) 的方法
mul	用于乘法 (*) 的方法

DataFrame和Series之间的运算 (Operations between DataFrame and Series)

DataFrame和Series之间算术运算也是有明确规定的。下面的例子是计算一个二维数组与其某行之间的差，这叫做广播 (broadcasting) 。

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

```
In [ ]: arr[0]
```

```
In [ ]: # 二维数组与其第一行之间的差, 这样的运算叫做广播 (broadcasting)
arr - arr[0]
```

- 在DataFrame 和 Series之间的运算也是类似的广播

```
In [ ]: frame = DataFrame(np.arange(12.).reshape((4, 3)), columns=list('bde'),  
index=['Utah', 'Ohio', 'Texas', 'Oregon'])  
frame
```

```
In [ ]: series = frame.ix[0] # 取出第一行  
series
```

```
In [ ]: # 减法运算也进行广播  
frame - series
```

- 默认情况下，DataFrame和Series之间的算术运算将会将Series 的索引匹配到DataFrame的列，然后沿着行一直向下广播。
- 如果某个索引值在DataFrame的列或Series的索引中找不到，则参与运算的两个对象就会被重新索引以形成并集。
- 如果希望匹配行并在列上广播，则必须使用算术运算方法。

```
In [ ]: # 如果某个索引值在DataFrame的列或Series的索引中找不到,  
# 则参与运算的两个对象就会被重新索引以形成并集。  
series2 = Series(range(3), index=['b', 'e', 'f'])  
frame + series2
```

```
In [ ]: frame - series2
```

```
In [ ]: # 如果希望匹配行并在列上广播, 则必须使用算术运算*方法*, 即, 不能直接 A - B。  
series3 = frame['d']  
frame
```

```
In [ ]: # sub 是减法的算术方法, 即frame-series3, 传入的轴号0就是希望匹配的轴,  
# 这里希望匹配frame的行索引并进行广播  
frame.sub(series3, axis=0)
```

函数应用和映射 (Function application and mapping)

- NumPy的ufunc (元素级数组方法) 也可以用于操作pandas对象。

```
In [ ]: frame = DataFrame(np.random.randn(4, 3), columns=list('bde'),  
index=['Utah', 'Ohio', 'Texas', 'Oregon'])  
frame
```

```
In [ ]: np.abs(frame)
```

- 另一个常用的操作是，将函数应用到各列或行所形成的一维数组上。DataFrame的apply方法即可实现此功能。
- 许多最为常见的数组统计功能都被实现成DataFrame的方法（如sum 和 mean），因此无需使用apply方法。

```
In [ ]: # 定义一个lambda 函数  
f = lambda x: x.max() - x.min()
```

```
In [ ]: # 默认函数是在DataFrame的各个列上apply, 即轴=0  
print frame, br  
frame.apply(f)
```

```
In [ ]: # 如果指定了应用的轴 (axis=1), 则按行应用  
frame.apply(f, axis=1)
```

- 除标量值外，传递给apply的函数还可以返回由多个值组成的Series。

```
In [ ]: def f(x):  
    return Series([x.min(), x.max()], index=['min', 'max'])  
frame.apply(f)
```

- 元素级的Python函数也是可以用的。想得到frame中各个浮点值的格式化字符串，使用applymap即可：

```
In [ ]: format = lambda x: '%.2f' % x  
frame.applymap(format)
```

- 之所以叫applymap，是因为Series有一个用于应用元素级函数的map方法：

```
In [ ]: frame['e'].map(format)
```

排序 (Sorting)

1. 按索引排序：可使用sort_index方法，它将返回一个已排序的新对象。对DataFrame，则可以根据任意一个轴上的索引进行排序（axis=0即行，axis=1即列）

2. 按值排序：可使用使用sort_values()方法（order方法过时）。任何缺失值默认都会被放到Series的末尾。对DataFrame上进行按值排序，如果希望根据一个或多个列中的值进行排序，将一个或多个列的名字传递给by选项即可。

对DataFrame，按行只能进行索引排序，按列可以进行索引排序，也可以按照指定列进行按值排序

```
In [ ]: s = Series(range(4), index=['d', 'a', 'b', 'c'])  
s.sort_index() # 对Series的行索引排序 (按照索引由小到大的顺序排序)
```

```
In [ ]: # 对于DataFrame, 则可以根据任意一个轴上的索引进行排序。默认是轴0  
frame = DataFrame(np.arange(8).reshape((2, 4)), index=['three', 'one'],  
                  columns=['d', 'a', 'b', 'c'])  
frame.sort_index() # 默认在轴0, 即行上进行行索引排序 (按照索引由小到大的顺序排序)  
frame.sort_index(axis=1) # 选择在轴1, 即列上进行列索引排序  
frame.sort_index(axis=1, ascending=False) # 数据默认是按升序排序, 也可以降序排序
```

按值排序

- 要按值对Series进行排序，使用sort_values()方法。（order方法过时）。排序时，任何缺失值默认都会被放到Series的末尾。
- 在DataFrame上进行按值排序，使用sort_values()方法。（sort_index方法过时），如果希望根据一个或多个列中的值进行排序，将一个或多个列的名字传递给by选项即可。

```
In [ ]: s = Series([4, 7, -3, 2])
s.sort_values() # sort_values按值排序, sort_index()是按索引排序
```

```
In [ ]: frame = DataFrame({'b': [4, 7, -3, 2], 'a': [0, 1, 0, 1]}, index=list('3012'))
print frame
frame.sort_index() # 默认在axis=0 即行轴上进行行索引排序
frame.sort_index(axis=1) # 在axis=1 即列轴上进行列索引排序
frame.sort_values(by='b') # by选项指定列, 按照该列进行按值排序
frame.sort_values(by=['a', 'b']) # 要按照多个列的值排序, 就把多个列的名字传递给by选项
frame.sort_values(by=['a', 'b'], ascending=False) # 按降序排列
```

排名 (Ranking)

- 排名 (ranking) 跟排序关系密切，且它会增设一个排名值（从1开始，直到数组中有效数据的数量）。跟 numpy.argsort 产生的间接排序索引差不多，但可根据某种规则破坏平级关系
- 默认情况下，rank是通过“为各组分配一个平均排名”的方式破坏平级关系的。如果有两个相同的值都排在第三位，则rank排名值会对它们的排名进行平均，即返回3.5
- 也可以根据值在原数据中出现的顺序给出排名（即，相同的值就按照它们出现的先后顺序而排序），method='first'
- 也可以按照降序进行排名，平级关系使用max即最大排名值

```
In [ ]: s = Series([7, -5, 7, 4, 2, 0, 4])
s.rank() # 对Series对象的rank排名(从1开始), 排名值对于两个平级的排名会进行平均
s.rank(method='first') # 也可以平级关系按照出现的先后顺序排序, method='first'
s.rank(ascending=False, method='max') # 也可以按照降序进行排名, 平级关系使用max即最大排名值
```

```
In [ ]: frame = DataFrame({'b': [4.3, 7, -3, 2], 'a': [0, 1, 0, 1],
                           'c': [-2, 5, 8, -2.5]})
#print frame
frame.rank() # 默认在每列上按照Series进行rank排名
frame.rank(axis=1) # 在轴1上排序, 即每行上按照Series进行rank排名
```

下表列出排名时用于破坏平级关系的method选项

method	说明
'average'	默认：在相等分组中，为各个值分配平均排名
'min'	使用整个分组的最小排名
'max'	使用整个分组的最大排名
'first'	按值在原始数据中的出现顺序分配排名

带有重复值的轴索引 (Axis indexes with duplicate values)

到目前为止，前面所有数据都有唯一的轴标签（索引值）。虽然很多pandas函数（如reindex）都要求标签唯一，但这并不是强制性的。

```
In [ ]: s = Series(range(5), index=['a', 'a', 'b', 'b', 'c']) # 带有重复索引值的Series
```

- `index.is_unique` 属性可以判断索引的值是否是唯一的
- 对于带有重复值的索引，数据选取的行为将会有些不同。如果某个索引对应多个值，则返回一个Series；而对应单个值，则返回一个标量值。

```
In [ ]: s.index.is_unique # index的is_unique 属性可以说明索引的值是否唯一
```

```
In [ ]: print s['a'], type(s['a']) # 如果某个索引对s应多个值，则返回一个Series  
print s['c'], type(s['c']) # 如果索引对应单个值，则返回一个标量值
```

- 对DataFrame的行进行索引时也是类似的。索引对应多个值，则返回DataFrame；如果对应单个值，则返回一个Series。

```
In [ ]: df = DataFrame(np.random.randn(4, 3), index=['a', 'a', 'b', 'c'])  
print df.ix['a'], type(df.ix['a']) # 索引对应多个值，则返回一个DataFrame  
print df.ix['c'], type(df.ix['c']) # 索引对应单个值，则返回一个Series
```

数据的汇总和描述统计

pandas对象拥有一组常用的数学和统计方法。它们大部分都属于约简和汇总统计，用于从Series中提取单个值（如 `sum` 和 `mean`）或从DataFrame的行或列中提取一个Series。跟对应的NumPy数组相比，它们都是基于没有缺失数据的假设而构建的。

- NA值会自动被排除，除非整个切片（行或列）都是NA。通过`skipna`选项可以禁用该功能。

```
In [ ]: df = DataFrame([[1.4, np.nan], [7.1, -4.5],  
                      [np.nan, np.nan], [0.75, -1.3]],  
                      index=['a', 'b', 'c', 'd'],  
                      columns=['one', 'two'])  
df
```

```
In [ ]: # 默认，一个列是一个series。调用DataFrame的sum方法将会返回一个含有列小计的Series:  
df.sum()
```

```
In [ ]: # 传入axis=1将会按行进行求和运算  
df.sum(axis=1)
```

- NA值会自动被排除，除非整个切片（行或列）都是NA。通过`skipna`选项可以禁用该功能。

```
In [ ]: # NA值会自动被排除，除非整个切片（行或列）都是NA。通过skipna选项可以禁用该功能。  
df.mean(axis=1, skipna=False)
```

下表列出了这些约简方法的常用选项。

选项	说明
axis	约简的轴。DataFrame的行用0，列用1
skipna	排除缺失值，默认值为True
level	如果轴是层次化索引的（即MultiIndex），则根据level分组约简

- 有些方法（如idxmin 和 idxmax）返回的是间接统计（比如达到最小值或最大值的索引）

```
In [ ]: # 返回的是间接统计，即，达到最大值的索引  
df.idxmax()
```

- 另一些方法是累计型的：

```
In [ ]: # 累计型方法  
df.cumsum()
```

- 还有一种方法，既不是约简型的，也不是累计型的。如 describe，用于一次性产生多个汇总统计的描述，返回值也是DataFrame型。
- 对于非数值型数据，describe会产生另外一种汇总统计。

```
In [ ]: # describe方法用于一次性产生多个汇总统计，返回的也是DataFrame型描述  
print type(df.describe()), br  
df.describe()
```

```
In [ ]: # 对于非数值型数据，describe会产生另外一种汇总统计  
obj = Series(['a', 'a', 'b', 'c'] * 4)  
# print obj, br  
obj.describe()
```

下表列出了所有与描述统计有关的方法

表列出描述和汇总统计

方法	说明
count	非NA值的数量
describe	针对Series或各DataFrame列计算汇总统计
min、max	计算最小值和最大值
argmin、argmax	计算能够获取到最小值和最大值的索引位置（整数）
idxmin、idxmax	计算能够获取到最小是和最大值的索引值
quantile	计算样本的分位数（0到1）
sum	值的总和
mean	值的平均数
median	值的算术中位数（50%分位数）
mad	根据平均值计算平均绝对离差
var	样本值的方差
std	样本值的标准差
skew	样本值的偏度（三阶矩）
kurt	样本值的峰度（四阶矩）
cumsucm	样本值的累计和
cummin、cummax	样本值的累计最大值和累计最小值
cumprod	样本值的累计积
diff	计算一阶差分（对时间序列很有用）
pct_change	计算百分数变化

相关系数与协方差（Correlation and covariance）

有些汇总统计（如相关系数和协方差）是通过参数对计算出来的。

- 下面的几个DataFrame的数据是来自Yahoo！Finance的股票价格和成交量。

采集股票数据

方法：使用Yahoo 金融提供的API

需要安装库 `install pandas_datareader`

- `import pandas_datareader as pdr`

或者

- `from pandas_datareader import data, wb`

```
In [ ]: # 方式： 使用Yahoo Finance的API获取四个公司的股票数据
import pandas as pd
import numpy as np
from pandas import Series, DataFrame

from pandas_datareader import data

all_stock = {}
for ticker in ['AAPL', 'IBM', 'MSFT', 'GOOG']:
    all_stock[ticker] = data.get_data_yahoo(ticker) # 默认从2010年1月起始, start='7/1/2005'

price = DataFrame({tic: data['Adj Close']
                  for tic, data in all_stock.iteritems()})
volume = DataFrame({tic: data['Volume']
                    for tic, data in all_stock.iteritems()})
open = DataFrame({tic: data['Open']
                  for tic, data in all_stock.iteritems()})
high = DataFrame({tic: data['High']
                  for tic, data in all_stock.iteritems()})
low = DataFrame({tic: data['Low']
                  for tic, data in all_stock.iteritems()})
```

```
In [ ]: # 最前面五条交易量数据
volume.head()
```

```
In [ ]: # 最后面五条价格数据
price.tail()
```

- 接下来计算价格的百分数变化

```
In [ ]: # 计算价格的百分数变化
returns = price.pct_change()
# 价格百分数变化的最后五条记录
returns.tail()
```

```
In [ ]: # 获得DataFrame的一个列MSFT, 得到一个Series
print type(returns.MSFT), br
# 获得MSFT最后五条数据, 返回一个Series
returns.MSFT.tail()
```

- Series的 `corr` 方法用于计算两个Series中重叠的、非NA的、按索引对齐的值的相关系数。
- 与此类似，`cov` 用于计算协方差。

```
In [ ]: # 计算微软与IBM之间的股票价格百分数变化的相关系数  
returns. MSFT. corr(returns. IBM)
```

```
In [ ]: # 计算微软与IBM之间的股票价格百分数变化的协方差  
returns. MSFT. cov(returns. IBM)
```

- 对于DataFrame的corr和cov方法将以DataFrame的形式返回完整的相关系数或协方差矩阵：

```
In [ ]: # 计算DataFrame的corr, 对角线为1.0  
returns. corr()
```

```
In [ ]: # 计算DataFrame的cov, 协方差矩阵  
returns. cov()
```

- 利用DataFrame的corrwith方法，可以计算其列或行跟另一个Series或DataFrame之间的相关系数。
- 传入一个Series将会返回一个相关系数值Series（针对各列进行计算）
- 传入一个DataFrame则会计算按列名配对的相关系数；如果传入axis=1，则计算按行的相关系数。

```
In [ ]: # 计算DataFrame与Series之间的相关系数  
returns. corrwith(returns. IBM)
```

```
In [ ]: volume. tail()
```

```
In [ ]: # 计算价格的百分数变化(DataFrame)与交易量(DataFrame)之间的相关系数, 产生按列名匹配的相关系数  
returns. corrwith(volume)
```

```
In [ ]: # 传入axis=1, 则计算按行的相关系数  
returns. corrwith(volume, axis=1). tail()
```

唯一值、值计数以及成员资格 (Unique values, value counts, and membership)

还有一类方法可以从一维Series的值中抽取信息，例如：

- 第一个函数是unique，可以得到Series中的唯一值数组numpy.ndarray类型，返回的唯一值是未排序的数组numpy.ndarray类型，如果需要，还可以对结果再次进行排序（unique.sort()）
- value_counts 用于计算一个Series中各值出现的频率。结果Series是按值频率降序排列的。
obj.value_counts()

value_counts 还是一个顶级pandas方法，可用于任何数组或序列。 pd.value_counts(obj.values, sort=False)

- 最后一个是isin，它用于判断矢量化集合的成员资格，可用于选取Series中或DataFrame列中数据的子集。

```
In [ ]: obj = Series(['c', 'a', 'd', 'a', 'a', 'b', 'b', 'c', 'c'])
```

```
In [ ]: # unique函数得到Series中的唯一值数组, 返回的唯一值是未排序的数组numpy.ndarray类型  
uniques = obj.unique()  
uniques
```

```
In [ ]: # value_counts函数用于计算一个Series中各值出现的频率，返回的结果Series是按值频率降序排列的
obj.value_counts()
```

```
In [ ]: # value_counts函数还是pandas的顶级方法，可以用于任何数组或序列
pd.value_counts(obj.values, sort=False)
```

```
In [ ]: # isin函数用于判断矢量化集合的成员资格，返回一个布尔型数组，用于选取Series或DataFrame列中数据的子集
mask = obj.isin(['b', 'c'])
mask
```

```
In [ ]: # 选取Series中数据的子集
obj[mask]
```

```
In [ ]: data = DataFrame({'Qu1': [1, 3, 4, 3, 4],
                       'Qu2': [2, 3, 1, 2, 3],
                       'Qu3': [1, 5, 2, 4, 4]})
```

```
In [ ]: # value_counts只能对Series计算，因此对于DataFrame需要指定在哪个列上进行值计数
data.Qu1.value_counts()
```

```
In [ ]: # value_counts也是pandas的顶级方法
pd.value_counts(data.Qu2, sort=True)
```

```
In [ ]: print data, br
result = data.apply(pd.value_counts).fillna(0)
result
```

下表列出这些方法的一些参考信息

唯一值、值计数、成员资格方法

方法	说明
isin	计算一个表示“Series各值是否包含于传入的值序列中”的布尔型数组
unique	计算Series中的唯一值数组，按发现的顺序返回
value_counts	返回一个Series，其索引为唯一值，其值为频率，按计数值降序排列

处理缺失值 (Handling missing data)

缺失数据 (missing data) 在大部分数据分析应用中都很常见。pandas的设计目标之一就是让缺失数据的处理任务尽量轻松。例如，pandas对象上的所有描述统计都排除了缺失数据。

- pandas使用浮点值 NaN (Not a Number) 表示浮点和非浮点数组中的缺失数据, np.nan, 是一个便于被检测出来的标记而已。
- Python内置的None值也会被当做NA处理, value = None
- pandas 的NA表现形式很简单也很可靠。由于NumPy的数据类型体系中缺乏真正的NA数据类型或位模式，所以目前最佳的解决方案可能就是pandas的NA (一套简单API以及足够全面的性能特征)。以后随着NumPy的发展，也许问题会变化。

```
In [ ]: string_data = Series(['aardvark', 'artichoke', np.nan, 'avocado'])  
string_data
```

```
In [ ]: string_data.isnull()
```

```
In [ ]: string_data[0] = None  
print string_data[0]  
string_data.isnull()
```

表列出了NA处理方法

方法	说明
dropna	根据各标签的值中是否存在缺失数据对轴标签进行过滤，可通过阈值调节对缺失值的容忍度
fillna	用指定值或插值方法 (如ffill或bfill) 填充缺失数据
isnull	返回一个含有布尔值的对象，这些布尔值表示哪些值是缺失值/NA，该对象的类型与源类型一样
notnull	isnull的否定式

滤除缺失数据 (Filtering out missing data)

过滤掉缺失数据的方法有很多种。

- 对于一个Series，dropna方法返回一个仅含非空数据和索引值的Series；也可以通过布尔型索引达到这个目的
- 但是对于一个DataFrame对象，事情有点复杂。dropna默认丢弃任何含有缺失值的行。

```
In [ ]: from numpy import nan as NA  
data = Series([1, NA, 3.5, NA, 7])  
# 对于一个Series, dropna方法返回一个仅含非空数据和索引值的Series  
data.dropna()
```

```
In [ ]: # 通过布尔型索引过滤掉空值  
#data.notnull()  
data[data.notnull()]
```

- 对于一个DataFrame对象，事情有点复杂。dropna默认丢弃任何含有缺失值的行；
- 传入how='all'将只丢弃全为NA的那些行；
- 如果要丢弃列，只需传入axis=1即可；
- 另一个滤除DataFrame行的问题涉及时间序列数据。如果只想留下部分观测数据，可以用thresh参数实现此目的。

```
In [ ]: # 对于一个DataFrame对象，事情有点复杂，dropna默认丢弃任何含有缺失值的行
         data = DataFrame([[1., 6.5, 3.], [1., NA, NA],
                           [NA, NA, NA], [NA, 6.5, 3.]])
         cleaned = data.dropna()
         data
```

```
In [ ]: # 丢弃任何含有缺失值的行
         cleaned
```

```
In [ ]: # 传入how='all'将只丢弃全为 NA 的那些行；how='any'是默认值
         data.dropna(how='all')
```

```
In [ ]: # 为data增加第四列
         data[4] = NA
         data
```

```
In [ ]: # 如果要丢弃列，只需传入axis=1即可
         data.dropna(axis=1, how='all')
```

```
In [ ]: df = DataFrame(np.random.randn(7, 3))
         df.ix[:4, 1] = NA; df.ix[:2, 2] = NA
         df
```

```
In [ ]: # 另一个滤除DataFrame行的问题涉及时间序列数据。如果只想留下部分观测数据，可以用thresh参数实现此目的
         # thresh=2 表示保留具有2个及以上非NA的行
         df.dropna(thresh=2)
```

填充缺失数据 (Filling in missing data)

我们可能不想滤除缺失数据（有可能会丢弃跟它有关的其他数据），而是希望通过其他方式填补那些“空洞”。

- 对于大多数情况而言，fillna方法是最主要的函数。通过一个常数调用fillna就会将缺失值替换为那个常数值，df.fillna(C)
- 若是通过一个字典调用fillna，就可以实现对不同的列填充不同的值
- fillna默认会返回新对象，但也可以对现有对象进行就地修改
- 对reindex有效的那些插值方法也可用于fillna
- 填充缺失值，也可以使用更聪明一些的办法，比如，可以传入Series的平均值或中位数，
data.fillna(data.median())

```
In [ ]: print df, br
         # 将缺失值替换为0
         df.fillna(0)
```

```
In [ ]: # 若是通过一个字典调用fillna，就可以实现对不同的列填充不同的值
         df.fillna({1: 0.5, 2: -1})
```

```
In [ ]: # fillna默认会返回新对象，但也可以对现有对象进行就地修改  
# always returns a reference to the filled object  
# 总是返回被填充对象的引用  
print df, br  
_ = df.fillna(0, inplace=True)  
df
```

```
In [ ]: df = DataFrame(np.random.randn(6, 3))  
df.ix[2:, 1] = NA; df.ix[4:, 2] = NA  
df
```

```
In [ ]: # 对reindex有效的那些插值方法也可用于fillna  
df.fillna(method='ffill')
```

```
In [ ]: # limit限定了前向或后向填充时可以连续填充的最大数量  
df.fillna(method='ffill', limit=2)
```

- 填充缺失值，也可以使用更聪明一些的办法，比如，可以传入Series的平均值或中位数

```
In [ ]: data = Series([1., NA, 3.5, NA, 7])  
# 使用平均值去填充缺失值  
data.fillna(data.mean())
```

```
In [ ]: # 使用中位数去填充缺失值  
data.fillna(data.median())
```

下表列出fillna函数的参数参考

表：fillna函数的参数

参数	说明
value	用于填充缺失值的标量值或字典对象
method	插值方式。如果函数调用时未指定其他参数的话，默认为“ffill”
axis	待填充的轴， 默认axis=0
inplace	修改调用者对象而不产生副本
limit	(对于前向和后向填充) 可以连续填充的最大数量

层次化索引 (Hierarchical indexing)

层次化索引 (hierarchical indexing) 是pandas的一项重要功能，它使你能做一个轴上拥有多个（两个以上）索引级别。抽象点说，它使你能以低维度形式处理高维度数据。

下面一个简单的例子，创建一个Series，并用一个由列表或数组组成的列表作为索引。

```
In [ ]: np.random.randn(3)
```

- 这就是带有MultiIndex索引的Series的格式化输出形式。索引之间的“间隔”表示“直接使用上面的标签”

```
In [ ]: data.index
```

- 对于一个层次化索引的对象，选取数据子集的操作很简单：

```
In [ ]: # 选择一行  
data['b']
```

```
In [ ]: # 选择几行的第一种方式  
data['b':'c']
```

```
In [ ]: # 选择几行的第二种方式  
data.ix[['b', 'd']]
```

```
In [ ]: # 按照“内层”另外一种索引选择数据  
data[:, 2]
```

- 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。比如下面这段数据可以通过其unstack方法被重新安排到一个DataFrame中。
 - unstack，即旋转（pivot），可以将带有MultiIndex索引的Series生成DataFrame，原有的外索引作为行索引保留，内层索引旋转为列索引（即轴发生旋转）。旋转后的层自动被排序。
 - unstack的逆运算是stack
 - 后面章节会继续详细讲解stack和unstack

```
In [ ]: #data.unstack?  
# Unstack, a.k.a. pivot, Series with MultiIndex to produce DataFrame.  
# The level involved will automatically get sorted.
```

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

In []: # unstack的逆运算
data.unstack().stack()

- 对于一个DataFrame，每条轴都可以有分层索引
 - 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（**不要将索引名称和轴标签混为一谈！**）

- 以上的是轴标签，每个轴没有名字(name) 不要将索引名称和轴标签混为一谈！

```
In [ ]: frame.index
```

```
In [ ]: frame.index.names = ['key1', 'key2']
frame.columns.names = ['state', 'color']
frame
```

- 有了分级的列索引，可以轻松选取列分组

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

- 可以单独创建MultiIndex然后复用。因此，前面的DataFrame中的分级列可以如下方式创建：

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

重排分级顺序 (Reordering and sorting levels)

有时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）
- sortlevel根据单个级别中的值对数据进行排序（稳定的）。交换级别时，常常也会用到sortlevel，这样最终结果就是有序的了。

```
In [ ]: print frame, br
frame.swaplevel('key1', 'key2')
# 在轴1上进行互换级别
# frame.swaplevel(0, 1, axis=1)
```

```
In [ ]: # sortlevel根据单个级别中的值对数据进行排序（稳定的）。
# 交换级别时，常常也会用到sortlevel，这样最终结果就是有序的了。
# 默认axis=0; level=1是指第二个索引层即key2作为primary key
frame.sortlevel(1)
# 在轴1上进行level=0的排序
# frame.sortlevel(axis=1, level=0)
```

```
In [ ]: # 交换第0层和第1层的索引，然后在第0层排序
frame.swaplevel(0, 1).sortlevel(0)
```

```
In [ ]: # 在轴1上交换第0层和第1层的索引，然后在轴1上的第0层排序
frame.swaplevel(0, 1, axis=1).sortlevel(0, axis=0)
```

根据级别汇总统计 (Summary statistics by level)

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。可以根据行或列上的级别来进行求和。

```
In [ ]: # 在轴0的level=1上求和
frame.sum(level='key2')
```

- 前面的操作只能一次指定在某个轴的一个级别上进行汇总
- 如果想在多个轴的多个级别上进行汇总，可以采用分步汇总，即先在某个指定轴的级别上汇总，得到一个新的DataFrame，然后针对新的DataFrame再次指定进行汇总的轴和级别。

```
In [ ]: # 在轴1的level=color上求和
f2=frame.sum(level='color', axis=1)
f2
```

```
In [ ]: f3=f2.sum(level='key2', axis=0)
f3
```

使用DataFrame的列 (Using a DataFrame's columns)

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。

- 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。
- 默认情况下，那些列会从DataFrame中移除，但也可以将其保留下来,drop=False。

```
In [ ]: frame = 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]})
frame
```

```
In [ ]: # 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。
frame2 = frame.set_index(['c', 'd'])
frame2
```

```
In [ ]: # 默认情况下，那些列会从DataFrame中移除，但也可以将其保留下来drop=False
frame.set_index(['c', 'd'], drop=False)
```

- reset_index的功能和set_index刚好相反，层次化索引的级别会被转移到列里面

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

其他有关pandas的话题 (Other pandas topics)

整数索引 (Integer indexing)

操作由整数索引的pandas对象常常会让新手抓狂，因为它们跟内置的Python数据结构（例如列表和元组）在索引语义上有些不同。

- 例如，下面这段代码虽然不会有错，在这种情况下，虽然pandas会求助于整数索引，因为这里我们有一个含有0、1、2的索引，但是pandas很难推断出用户想要什么（基于标签或位置的索引）

```
In [ ]: #
ser = Series(np.arange(3.))
ser.iloc[-1]
```

```
In [ ]: ser
```

- 相反，对于一个非整数索引，就没有这样的歧义：

```
In [ ]: ser2 = Series(np.arange(3.), index=['a', 'b', 'c'])  
ser2[-1]
```

- 为了保持良好的一致性，如果你的轴索引含有索引器，那么根据整数进行数据选取的操作将总是面向标签的。这也包括用ix切片：

```
In [ ]: ser.ix[:1]
```

```
In [ ]: ser.ix[1:2]
```

- 如果需要可靠的、不必考虑索引类型的、基于位置的索引，可以采用如下的方法。
- 对于Series，可以使用（1）`iget_value()`方法和（2）`iloc[]`方法，两者等同
- 对于DataFrame的`irow`和`icol`方法：

```
In [ ]: ser3 = Series(range(3), index=[-5, 1, 3])  
ser3  
#下面两个方法等同
```

```
In [ ]: ser3.iget_value(2)
```

```
In [ ]: ser3.iloc[2]
```

```
ser3.iget_value(i, axis=0) # 返回轴0上的第i-1个位置的值，索引从0开始 # Return the i-th value or values in the Series  
by location  
ser3.iloc[2] # 返回整数索引位置的值 # Purely integer-location based indexing for selection by position.
```

```
In [ ]: frame = DataFrame(np.arange(6).reshape((3, 2)), index=[2, 0, 1])  
frame
```

```
In [ ]: # 默认返回第一行的值  
frame.iloc[0]
```

```
In [ ]: # 返回第i+1行的值  
frame.irow(0)
```

```
In [ ]: # 返回第i+1列的值  
frame.icol(1)
```

面板数据 (Panel data)

- pandas有一个Panel数据结构，可以将其看做一个三维版的DataFrame。
- pandas的大部分开发工作都集中在表格型数据的操作上，因为这些数据更常见，而且层次化索引也使得多数情况下没有必要使用真正的N维数组。
- 可以用一个由DataFrame对象组成的字典或一个三维ndarray来创建Panel对象
- Panel中的每一项（类似于DataFrame的列）都是一个DataFrame

```
In [ ]: import pandas as pd
from pandas_datareader import data

# 用一个由DataFrame对象组成的字典或一个三维ndarray来创建Panel对象
pdata = pd.Panel(dict((stk, data.get_data_yahoo(stk))
                      for stk in ['AAPL', 'GOOG', 'MSFT', 'DELL']))
```

```
In [ ]: # Panel中的每一项（类似于DataFrame的列）都是一个DataFrame
pdata
```

```
In [ ]: # 交换轴
pdata = pdata.swapaxes('items', 'minor')
pdata['Adj Close']
```

- 基于ix的标签索引被推广到了三个维度，因此我们可以选取指定日期或日期范围的所有数据：

```
In [ ]: pdata.ix[:, '6/1/2012', :]
```

```
In [ ]: pdata.ix['Adj Close', '5/22/2012':, :]
```

- 另一个用于呈现面板数据（尤其是对逆和统计模型）的办法是“堆积式的”DataFrame形式：

```
In [ ]: stacked = pdata.ix[:, '5/30/2012':, :].to_frame()
stacked
```

- DataFrame有一个相应的to_panel方法，它是to_frame的逆运算

```
In [ ]: stacked.to_panel()
```