

*Agile Tools for Real-World Data*

# Python for Data Analysis



O'REILLY®

*Wes McKinney*

# Getting Started with pandas

pandas will be the primary library of interest throughout much of the rest of the book. It contains high-level data structures and manipulation tools designed to make data analysis fast and easy in Python. pandas is built on top of NumPy and makes it easy to use in NumPy-centric applications.

As a bit of background, I started building pandas in early 2008 during my tenure at AQR, a quantitative investment management firm. At the time, I had a distinct set of requirements that were not well-addressed by any single tool at my disposal:

- Data structures with labeled axes supporting automatic or explicit data alignment. This prevents common errors resulting from misaligned data and working with differently-indexed data coming from different sources.
- Integrated time series functionality.
- The same data structures handle both time series data and non-time series data.
- Arithmetic operations and reductions (like summing across an axis) would pass on the metadata (axis labels).
- Flexible handling of missing data.
- Merge and other relational operations found in popular database databases (SQL-based, for example).

I wanted to be able to do all of these things in one place, preferably in a language well-suited to general purpose software development. Python was a good candidate language for this, but at that time there was not an integrated set of data structures and tools providing this functionality.

Over the last four years, pandas has matured into a quite large library capable of solving a much broader set of data handling problems than I ever anticipated, but it has expanded in its scope without compromising the simplicity and ease-of-use that I desired from the very beginning. I hope that after reading this book, you will find it to be just as much of an indispensable tool as I do.

Throughout the rest of the book, I use the following import conventions for pandas:

```
In [1]: from pandas import Series, DataFrame
```

```
In [2]: import pandas as pd
```

Thus, whenever you see `pd.` in code, it's referring to pandas. Series and DataFrame are used so much that I find it easier to import them into the local namespace.

## Introduction to pandas Data Structures

To get started with pandas, you will need to get comfortable with its two workhorse data structures: *Series* and *DataFrame*. While they are not a universal solution for every problem, they provide a solid, easy-to-use basis for most applications.

### Series

A Series is a one-dimensional array-like object containing an array of data (of any NumPy data type) and an associated array of data labels, called its *index*. The simplest Series is formed from only an array of data:

```
In [4]: obj = Series([4, 7, -5, 3])
```

```
In [5]: obj
```

```
Out[5]:
```

```
0    4  
1    7  
2   -5  
3    3
```

The string representation of a Series displayed interactively shows the index on the left and the values on the right. Since we did not specify an index for the data, a default one consisting of the integers 0 through N - 1 (where N is the length of the data) is created. You can get the array representation and index object of the Series via its `values` and `index` attributes, respectively:

```
In [6]: obj.values
```

```
Out[6]: array([ 4,  7, -5,  3])
```

```
In [7]: obj.index
```

```
Out[7]: Int64Index([0, 1, 2, 3])
```

Often it will be desirable to create a Series with an index identifying each data point:

```
In [8]: obj2 = Series([4, 7, -5, 3], index=['d', 'b', 'a', 'c'])
```

```
In [9]: obj2
```

```
Out[9]:
```

```
d    4  
b    7  
a   -5  
c    3
```

```
In [10]: obj2.index  
Out[10]: Index([d, b, a, c], dtype=object)
```

Compared with a regular NumPy array, you can use values in the index when selecting single values or a set of values:

```
In [11]: obj2['a']  
Out[11]: -5
```

```
In [12]: obj2['d'] = 6
```

```
In [13]: obj2[['c', 'a', 'd']]  
Out[13]:  
c    3  
a   -5  
d    6
```

NumPy array operations, such as filtering with a boolean array, scalar multiplication, or applying math functions, will preserve the index-value link:

```
In [14]: obj2  
Out[14]:  
d    6  
b    7  
a   -5  
c    3
```

```
In [15]: obj2[obj2 > 0]  
Out[15]:  
d    6  
b    7  
c    3
```

```
In [16]: obj2 * 2  
Out[16]:  
d    12  
b    14  
a   -10  
c     6
```

```
In [17]: np.exp(obj2)  
Out[17]:  
d    403.428793  
b  1096.633158  
a    0.006738  
c   20.085537
```

Another way to think about a Series is as a fixed-length, ordered dict, as it is a mapping of index values to data values. It can be substituted into many functions that expect a dict:

```
In [18]: 'b' in obj2  
Out[18]: True
```

```
In [19]: 'e' in obj2  
Out[19]: False
```

Should you have data contained in a Python dict, you can create a Series from it by passing the dict:

```
In [20]: sdata = {'Ohio': 35000, 'Texas': 71000, 'Oregon': 16000, 'Utah': 5000}
```

```
In [21]: obj3 = Series(sdata)
```

```
In [22]: obj3  
Out[22]:  
Ohio      35000  
Oregon    16000
```

```
Texas      71000  
Utah       5000
```

When only passing a dict, the index in the resulting Series will have the dict's keys in sorted order.

```
In [23]: states = ['California', 'Ohio', 'Oregon', 'Texas']
```

```
In [24]: obj4 = Series(sdata, index=states)
```

```
In [25]: obj4  
Out[25]:  
California      NaN  
Ohio            35000  
Oregon          16000  
Texas           71000
```

In this case, 3 values found in `sdata` were placed in the appropriate locations, but since no value for 'California' was found, it appears as `NaN` (not a number) which is considered in pandas to mark missing or `NA` values. I will use the terms "missing" or "`NA`" to refer to missing data. The `isnull` and `notnull` functions in pandas should be used to detect missing data:

```
In [26]: pd.isnull(obj4)      In [27]: pd.notnull(obj4)  
Out[26]:  
California    True           California   False  
Ohio          False          Ohio         True  
Oregon        False          Oregon       True  
Texas         False          Texas        True
```

Series also has these as instance methods:

```
In [28]: obj4.isnull()  
Out[28]:  
California    True  
Ohio          False  
Oregon        False  
Texas         False
```

I discuss working with missing data in more detail later in this chapter.

A critical Series feature for many applications is that it automatically aligns differently-indexed data in arithmetic operations:

```
In [29]: obj3      In [30]: obj4  
Out[29]:  
Ohio      35000    California     NaN  
Oregon    16000    Ohio          35000  
Texas     71000    Oregon         16000  
Utah      5000     Texas          71000
```

```
In [31]: obj3 + obj4  
Out[31]:  
California      NaN  
Ohio            70000  
Oregon          32000
```

```
Texas      142000  
Utah       NaN
```

Data alignment features are addressed as a separate topic.

Both the Series object itself and its index have a `name` attribute, which integrates with other key areas of pandas functionality:

```
In [32]: obj4.name = 'population'  
  
In [33]: obj4.index.name = 'state'  
  
In [34]: obj4  
Out[34]:  
state  
California      NaN  
Ohio            35000  
Oregon          16000  
Texas           71000  
Name: population
```

A Series's index can be altered in place by assignment:

```
In [35]: obj.index = ['Bob', 'Steve', 'Jeff', 'Ryan']  
  
In [36]: obj  
Out[36]:  
Bob     4  
Steve   7  
Jeff    -5  
Ryan    3
```

## DataFrame

A DataFrame represents a tabular, spreadsheet-like data structure containing an ordered collection of columns, each of which can be a different value type (numeric, string, boolean, etc.). The DataFrame has both a row and column index; it can be thought of as a dict of Series (one for all sharing the same index). Compared with other such DataFrame-like structures you may have used before (like R's `data.frame`), row-oriented and column-oriented operations in DataFrame are treated roughly symmetrically. Under the hood, the data is stored as one or more two-dimensional blocks rather than a list, dict, or some other collection of one-dimensional arrays. The exact details of DataFrame's internals are far outside the scope of this book.



While DataFrame stores the data internally in a two-dimensional format, you can easily represent much higher-dimensional data in a tabular format using hierarchical indexing, a subject of a later section and a key ingredient in many of the more advanced data-handling features in pandas.

There are numerous ways to construct a DataFrame, though one of the most common is from a dict of equal-length lists or NumPy arrays

```
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)
```

The resulting DataFrame will have its index assigned automatically as with Series, and the columns are placed in sorted order:

```
In [38]: frame
Out[38]:
   pop  state  year
0  1.5    Ohio  2000
1  1.7    Ohio  2001
2  3.6    Ohio  2002
3  2.4  Nevada  2001
4  2.9  Nevada  2002
```

If you specify a sequence of columns, the DataFrame's columns will be exactly what you pass:

```
In [39]: DataFrame(data, columns=['year', 'state', 'pop'])
Out[39]:
   year  state  pop
0  2000    Ohio  1.5
1  2001    Ohio  1.7
2  2002    Ohio  3.6
3  2001  Nevada  2.4
4  2002  Nevada  2.9
```

As with Series, if you pass a column that isn't contained in `data`, it will appear with NA values in the result:

```
In [40]: frame2 = DataFrame(data, columns=['year', 'state', 'pop', 'debt'],
                           ....:                 index=['one', 'two', 'three', 'four', 'five'])
In [41]: frame2
Out[41]:
   year  state  pop  debt
one    2000    Ohio  1.5    NaN
two    2001    Ohio  1.7    NaN
three  2002    Ohio  3.6    NaN
four   2001  Nevada  2.4    NaN
five   2002  Nevada  2.9    NaN
```

```
In [42]: frame2.columns
Out[42]: Index(['year', 'state', 'pop', 'debt'], dtype=object)
```

A column in a DataFrame can be retrieved as a Series either by dict-like notation or by attribute:

In [43]: frame2['state'] Out[43]: one      Ohio	In [44]: frame2.year Out[44]: one      2000
---	---

```
two      Ohio          two    2001
three    Ohio          three   2002
four     Nevada        four    2001
five     Nevada        five    2002
Name: state           Name: year
```

Note that the returned Series have the same index as the DataFrame, and their `name` attribute has been appropriately set.

Rows can also be retrieved by position or name by a couple of methods, such as the `ix` indexing field (much more on this later):

```
In [45]: frame2.ix['three']
Out[45]:
year    2002
state   Ohio
pop     3.6
debt    NaN
Name: three
```

Columns can be modified by assignment. For example, the empty '`debt`' column could be assigned a scalar value or an array of values:

```
In [46]: frame2['debt'] = 16.5

In [47]: frame2
Out[47]:
       year  state  pop  debt
one    2000   Ohio  1.5  16.5
two    2001   Ohio  1.7  16.5
three  2002   Ohio  3.6  16.5
four   2001  Nevada 2.4  16.5
five   2002  Nevada 2.9  16.5

In [48]: frame2['debt'] = np.arange(5.)

In [49]: frame2
Out[49]:
       year  state  pop  debt
one    2000   Ohio  1.5    0
two    2001   Ohio  1.7    1
three  2002   Ohio  3.6    2
four   2001  Nevada 2.4    3
five   2002  Nevada 2.9    4
```

When assigning lists or arrays to a column, the value's length must match the length of the DataFrame. If you assign a Series, it will be instead conformed exactly to the DataFrame's index, inserting missing values in any holes:

```
In [50]: val = Series([-1.2, -1.5, -1.7], index=['two', 'four', 'five'])

In [51]: frame2['debt'] = val

In [52]: frame2
Out[52]:
       year  state  pop  debt
two    -1.2   Ohio  1.5  -1.2
four   -1.5  Nevada 2.4  -1.5
five  -1.7  Nevada 2.9  -1.7
```

```
one    2000    Ohio  1.5   NaN
two    2001    Ohio  1.7  -1.2
three  2002    Ohio  3.6   NaN
four   2001  Nevada  2.4  -1.5
five   2002  Nevada  2.9  -1.7
```

Assigning a column that doesn't exist will create a new column. The `del` keyword will delete columns as with a dict:

```
In [53]: frame2['eastern'] = frame2.state == 'Ohio'
```

```
In [54]: frame2
```

```
Out[54]:
```

	year	state	pop	debt	eastern
one	2000	Ohio	1.5	NaN	True
two	2001	Ohio	1.7	-1.2	True
three	2002	Ohio	3.6	NaN	True
four	2001	Nevada	2.4	-1.5	False
five	2002	Nevada	2.9	-1.7	False

```
In [55]: del frame2['eastern']
```

```
In [56]: frame2.columns
```

```
Out[56]: Index(['year', 'state', 'pop', 'debt'], dtype=object)
```



The column returned when indexing a DataFrame is a *view* on the underlying data, not a copy. Thus, any in-place modifications to the Series will be reflected in the DataFrame. The column can be explicitly copied using the Series's `copy` method.

Another common form of data is a nested dict of dicts format:

```
In [57]: pop = {'Nevada': {2001: 2.4, 2002: 2.9},
....:           'Ohio': {2000: 1.5, 2001: 1.7, 2002: 3.6}}
```

If passed to DataFrame, it will interpret the outer dict keys as the columns and the inner keys as the row indices:

```
In [58]: frame3 = DataFrame(pop)
```

```
In [59]: frame3
```

```
Out[59]:
```

	Nevada	Ohio
2000	NaN	1.5
2001	2.4	1.7
2002	2.9	3.6

Of course you can always transpose the result:

```
In [60]: frame3.T
Out[60]:
```

	2000	2001	2002
Nevada	NaN	2.4	2.9
Ohio	1.5	1.7	3.6

The keys in the inner dicts are unioned and sorted to form the index in the result. This isn't true if an explicit index is specified:

```
In [61]: DataFrame(pop, index=[2001, 2002, 2003])
Out[61]:
      Nevada  Ohio
2001      2.4   1.7
2002      2.9   3.6
2003      NaN   NaN
```

Dicts of Series are treated much in the same way:

```
In [62]: pdata = {'Ohio': frame3['Ohio'][:-1],
.....:           'Nevada': frame3['Nevada'][:2]}

In [63]: DataFrame(pdata)
Out[63]:
      Nevada  Ohio
2000      NaN   1.5
2001      2.4   1.7
```

For a complete list of things you can pass the DataFrame constructor, see [Table 5-1](#).

If a DataFrame's `index` and `columns` have their `name` attributes set, these will also be displayed:

```
In [64]: frame3.index.name = 'year'; frame3.columns.name = 'state'

In [65]: frame3
Out[65]:
      state  Nevada  Ohio
year
2000      NaN   1.5
2001      2.4   1.7
2002      2.9   3.6
```

Like Series, the `values` attribute returns the data contained in the DataFrame as a 2D ndarray:

```
In [66]: frame3.values
Out[66]:
array([[ nan,  1.5],
       [ 2.4,  1.7],
       [ 2.9,  3.6]])
```

If the DataFrame's columns are different dtypes, the dtype of the values array will be chosen to accomodate all of the columns:

```
In [67]: frame2.values
Out[67]:
array([[2000, Ohio, 1.5, nan],
       [2001, Ohio, 1.7, -1.2],
       [2002, Ohio, 3.6, nan],
       [2001, Nevada, 2.4, -1.5],
       [2002, Nevada, 2.9, -1.7]], dtype=object)
```

Table 5-1. Possible data inputs to DataFrame constructor

Type	Notes
2D ndarray	A matrix of data, passing optional row and column labels
dict of arrays, lists, or tuples	Each sequence becomes a column in the DataFrame. All sequences must be the same length.
NumPy structured/record array	Treated as the “dict of arrays” case
dict of Series	Each value becomes a column. Indexes from each Series are unioned together to form the result’s row index if no explicit index is passed.
dict of dicts	Each inner dict becomes a column. Keys are unioned to form the row index as in the “dict of Series” case.
list of dicts or Series	Each item becomes a row in the DataFrame. Union of dict keys or Series indexes become the DataFrame’s column labels
List of lists or tuples	Treated as the “2D ndarray” case
Another DataFrame	The DataFrame’s indexes are used unless different ones are passed
NumPy MaskedArray	Like the “2D ndarray” case except masked values become NA/missing in the DataFrame result

## Index Objects

pandas’s Index objects are responsible for holding the axis labels and other metadata (like the axis name or names). Any array or other sequence of labels used when constructing a Series or DataFrame is internally converted to an Index:

```
In [68]: obj = Series(range(3), index=['a', 'b', 'c'])

In [69]: index = obj.index

In [70]: index
Out[70]: Index([a, b, c], dtype=object)

In [71]: index[1:]
Out[71]: Index([b, c], dtype=object)
```

Index objects are immutable and thus can’t be modified by the user:

```
In [72]: index[1] = 'd'
-----
Exception Traceback (most recent call last)
<ipython-input-72-676fdeb26a68> in <module>()
----> 1 index[1] = 'd'
/Users/wesm/code/pandas/pandas/core/index.pyc in __setitem__(self, key, value)
    302     def __setitem__(self, key, value):
    303         """Disable the setting of values."""
--> 304         raise Exception(str(self.__class__) + ' object is immutable')
    305
    306     def __getitem__(self, key):
Exception: <class 'pandas.core.index.Index'> object is immutable
```

Immutability is important so that Index objects can be safely shared among data structures:

```
In [73]: index = pd.Index(np.arange(3))

In [74]: obj2 = Series([1.5, -2.5, 0], index=index)

In [75]: obj2.index is index
Out[75]: True
```

**Table 5-2** has a list of built-in Index classes in the library. With some development effort, Index can even be subclassed to implement specialized axis indexing functionality.



Many users will not need to know much about Index objects, but they're nonetheless an important part of pandas's data model.

*Table 5-2. Main Index objects in pandas*

Class	Description
Index	The most general Index object, representing axis labels in a NumPy array of Python objects.
Int64Index	Specialized Index for integer values.
MultiIndex	"Hierarchical" index object representing multiple levels of indexing on a single axis. Can be thought of as similar to an array of tuples.
DatetimeIndex	Stores nanosecond timestamps (represented using NumPy's datetime64 dtype).
PeriodIndex	Specialized Index for Period data (timespans).

In addition to being array-like, an Index also functions as a fixed-size set:

```
In [76]: frame3
Out[76]:
state  Nevada  Ohio
year
2000      NaN   1.5
2001      2.4   1.7
2002      2.9   3.6

In [77]: 'Ohio' in frame3.columns
Out[77]: True

In [78]: 2003 in frame3.index
Out[78]: False
```

Each Index has a number of methods and properties for set logic and answering other common questions about the data it contains. These are summarized in [Table 5-3](#).

Table 5-3. Index methods and properties

Method	Description
append	Concatenate with additional Index objects, producing a new Index
diff	Compute set difference as an Index
intersection	Compute set intersection
union	Compute set union
isin	Compute boolean array indicating whether each value is contained in the passed collection
delete	Compute new Index with element at index <i>i</i> deleted
drop	Compute new index by deleting passed values
insert	Compute new Index by inserting element at index <i>i</i>
is_monotonic	Returns True if each element is greater than or equal to the previous element
is_unique	Returns True if the Index has no duplicate values
unique	Compute the array of unique values in the Index

## Essential Functionality

In this section, I'll walk you through the fundamental mechanics of interacting with the data contained in a Series or DataFrame. Upcoming chapters will delve more deeply into data analysis and manipulation topics using pandas. This book is not intended to serve as exhaustive documentation for the pandas library; I instead focus on the most important features, leaving the less common (that is, more esoteric) things for you to explore on your own.

### Reindexing

A critical method on pandas objects is `reindex`, which means to create a new object with the data *conformed* to a new index. Consider a simple example from above:

```
In [79]: obj = Series([4.5, 7.2, -5.3, 3.6], index=['d', 'b', 'a', 'c'])

In [80]: obj
Out[80]:
d    4.5
b    7.2
a   -5.3
c    3.6
```

Calling `reindex` on this Series rearranges the data according to the new index, introducing missing values if any index values were not already present:

```
In [81]: obj2 = obj.reindex(['a', 'b', 'c', 'd', 'e'])

In [82]: obj2
Out[82]:
a   -5.3
```

```

b    7.2
c    3.6
d    4.5
e    NaN

In [83]: obj.reindex(['a', 'b', 'c', 'd', 'e'], fill_value=0)
Out[83]:
a   -5.3
b    7.2
c    3.6
d    4.5
e    0.0

```

For ordered data like time series, it may be desirable to do some interpolation or filling of values when reindexing. The `method` option allows us to do this, using a method such as `ffill` which forward fills the values:

```

In [84]: obj3 = Series(['blue', 'purple', 'yellow'], index=[0, 2, 4])

In [85]: obj3.reindex(range(6), method='ffill')
Out[85]:
0    blue
1    blue
2  purple
3  purple
4  yellow
5  yellow

```

**Table 5-4.** lists available `method` options. At this time, interpolation more sophisticated than forward- and backfilling would need to be applied after the fact.

*Table 5-4. reindex method (interpolation) options*

Argument	Description
<code>ffill</code> or <code>pad</code>	Fill (or carry) values forward
<code>bfill</code> or <code>backfill</code>	Fill (or carry) values backward

With DataFrame, `reindex` can alter either the (row) index, columns, or both. When passed just a sequence, the rows are reindexed in the result:

```

In [86]: frame = DataFrame(np.arange(9).reshape((3, 3)), index=['a', 'c', 'd'],
....:                           columns=['Ohio', 'Texas', 'California'])

In [87]: frame
Out[87]:
   Ohio  Texas  California
a      0      1      2
c      3      4      5
d      6      7      8

In [88]: frame2 = frame.reindex(['a', 'b', 'c', 'd'])

In [89]: frame2
Out[89]:

```

```

      Ohio  Texas  California
a      0      1          2
b     NaN    NaN        NaN
c      3      4          5
d      6      7          8

```

The columns can be reindexed using the `columns` keyword:

```
In [90]: states = ['Texas', 'Utah', 'California']
```

```
In [91]: frame.reindex(columns=states)
```

```
Out[91]:
```

```

      Texas  Utah  California
a      1  NaN      2
c      4  NaN      5
d      7  NaN      8

```

Both can be reindexed in one shot, though interpolation will only apply row-wise (axis 0):

```
In [92]: frame.reindex(index=['a', 'b', 'c', 'd'], method='ffill',
....:                 columns=states)
```

```
Out[92]:
```

```

      Texas  Utah  California
a      1  NaN      2
b      1  NaN      2
c      4  NaN      5
d      7  NaN      8

```

As you'll see soon, reindexing can be done more succinctly by label-indexing with `ix`:

```
In [93]: frame.ix[['a', 'b', 'c', 'd'], states]
```

```
Out[93]:
```

```

      Texas  Utah  California
a      1  NaN      2
b     NaN    NaN    NaN
c      4  NaN      5
d      7  NaN      8

```

*Table 5-5. reindex function arguments*

Argument	Description
<code>index</code>	New sequence to use as index. Can be <code>Index</code> instance or any other sequence-like Python data structure. An <code>Index</code> will be used exactly as is without any copying.
<code>method</code>	Interpolation (fill) method, see <a href="#">Table 5-4</a> for options.
<code>fill_value</code>	Substitute value to use when introducing missing data by reindexing
<code>limit</code>	When forward- or backfilling, maximum size gap to fill
<code>level</code>	Match simple <code>Index</code> on level of <code>MultIndex</code> , otherwise select subset of
<code>copy</code>	Do not copy underlying data if new index is equivalent to old index. <code>True</code> by default (i.e. always copy data).

## Dropping entries from an axis

Dropping one or more entries from an axis is easy if you have an index array or list without those entries. As that can require a bit of munging and set logic, the `drop` method will return a new object with the indicated value or values deleted from an axis:

```
In [94]: obj = Series(np.arange(5.), index=['a', 'b', 'c', 'd', 'e'])
```

```
In [95]: new_obj = obj.drop('c')
```

```
In [96]: new_obj
```

```
Out[96]:
```

```
a    0  
b    1  
d    3  
e    4
```

```
In [97]: obj.drop(['d', 'c'])
```

```
Out[97]:
```

```
a    0  
b    1  
e    4
```

With DataFrame, index values can be deleted from either axis:

```
In [98]: data = DataFrame(np.arange(16).reshape((4, 4)),  
....:                  index=['Ohio', 'Colorado', 'Utah', 'New York'],  
....:                  columns=['one', 'two', 'three', 'four'])
```

```
In [99]: data.drop(['Colorado', 'Ohio'])
```

```
Out[99]:
```

	one	two	three	four
Utah	8	9	10	11
New York	12	13	14	15

```
In [100]: data.drop('two', axis=1)
```

```
Out[100]:
```

	one	three	four
Ohio	0	2	3
Colorado	4	6	7
Utah	8	10	11
New York	12	14	15

```
In [101]: data.drop(['two', 'four'], axis=1)
```

```
Out[101]:
```

	one	three
Ohio	0	2
Colorado	4	6
Utah	8	10
New York	12	14

## Indexing, selection, and filtering

Series indexing (`obj[...]`) works analogously to NumPy array indexing, except you can use the Series's index values instead of only integers. Here are some examples this:

```
In [102]: obj = Series(np.arange(4.), index=['a', 'b', 'c', 'd'])
```

```
In [103]: obj['b']
```

```
Out[103]: 1.0
```

```
In [104]: obj[1]
```

```
Out[104]: 1.0
```

```
In [105]: obj[2:4]
```

```
Out[105]:
```

```
In [106]: obj[['b', 'a', 'd']]
```

```
Out[106]:
```

```

c    2          b    1
d    3          a    0
                  d    3

In [107]: obj[[1, 3]]      In [108]: obj[obj < 2]
Out[107]:                   Out[108]:
b    1          a    0
d    3          b    1

```

Slicing with labels behaves differently than normal Python slicing in that the endpoint is inclusive:

```

In [109]: obj['b':'c']
Out[109]:
b    1
c    2

```

*Setting* using these methods works just as you would expect:

```

In [110]: obj['b':'c'] = 5

In [111]: obj
Out[111]:
a    0
b    5
c    5
d    3

```

As you've seen above, indexing into a DataFrame is for retrieving one or more columns either with a single value or sequence:

```

In [112]: data = DataFrame(np.arange(16).reshape((4, 4)),
.....:                      index=['Ohio', 'Colorado', 'Utah', 'New York'],
.....:                      columns=['one', 'two', 'three', 'four'])

In [113]: data
Out[113]:
   one  two  three  four
Ohio     0    1     2     3
Colorado  4    5     6     7
Utah     8    9    10    11
New York 12   13    14    15

In [114]: data['two']           In [115]: data[['three', 'one']]
Out[114]:                      Out[115]:
Ohio      1                      three  one
Colorado  5                      Ohio      2    0
Utah     9                      Colorado  6    4
New York 13                     Utah     10   8
Name: two                         New York 14   12

```

Indexing like this has a few special cases. First selecting rows by slicing or a boolean array:

```

In [116]: data[:2]           In [117]: data[data['three'] > 5]
Out[116]:                      Out[117]:
   one  two  three  four          one  two  three  four

```

Ohio	0	1	2	3	Colorado	4	5	6	7
Colorado	4	5	6	7	Utah	8	9	10	11
					New York	12	13	14	15

This might seem inconsistent to some readers, but this syntax arose out of practicality and nothing more. Another use case is in indexing with a boolean DataFrame, such as one produced by a scalar comparison:

```
In [118]: data < 5
Out[118]:
      one   two  three  four
Ohio    True  True  True  True
Colorado  True False False False
Utah    False False False False
New York False False False False
```

```
In [119]: data[data < 5] = 0
```

```
In [120]: data
Out[120]:
      one   two  three  four
Ohio     0     0     0     0
Colorado  0     5     6     7
Utah     8     9    10    11
New York 12    13    14    15
```

This is intended to make DataFrame syntactically more like an ndarray in this case.

For DataFrame label-indexing on the rows, I introduce the special indexing field `ix`. It enables you to select a subset of the rows and columns from a DataFrame with NumPy-like notation plus axis labels. As I mentioned earlier, this is also a less verbose way to do reindexing:

```
In [121]: data.ix['Colorado', ['two', 'three']]
Out[121]:
two    5
three  6
Name: Colorado
```

```
In [122]: data.ix[['Colorado', 'Utah'], [3, 0, 1]]
Out[122]:
      four  one  two
Colorado    7    0    5
Utah       11   8    9
```

```
In [123]: data.ix[2]
Out[123]:
one    8
two    9
three  10
four   11
Name: Utah
```

```
In [124]: data.ix[:'Utah', 'two']
Out[124]:
Ohio        0
Colorado    5
Utah       9
Name: two
```

```
In [125]: data.ix[data.three > 5, :3]
Out[125]:
```

```

      one  two  three
Colorado    0    5     6
Utah        8    9    10
New York   12   13    14

```

So there are many ways to select and rearrange the data contained in a pandas object. For DataFrame, there is a short summary of many of them in [Table 5-6](#). You have a number of additional options when working with hierarchical indexes as you'll later see.



When designing pandas, I felt that having to type `frame[:, col]` to select a column was too verbose (and error-prone), since column selection is one of the most common operations. Thus I made the design trade-off to push all of the rich label-indexing into `ix`.

*Table 5-6. Indexing options with DataFrame*

Type	Notes
<code>obj[val]</code>	Select single column or sequence of columns from the DataFrame. Special case conveniences: boolean array (filter rows), slice (slice rows), or boolean DataFrame (set values based on some criterion).
<code>obj.ix[val]</code>	Selects single row of subset of rows from the DataFrame.
<code>obj.ix[:, val]</code>	Selects single column of subset of columns.
<code>obj.ix[val1, val2]</code>	Select both rows and columns.
<code>reindex</code> method	Conform one or more axes to new indexes.
<code>xs</code> method	Select single row or column as a Series by label.
<code>icol, irow</code> methods	Select single column or row, respectively, as a Series by integer location.
<code>get_value, set_value</code> methods	Select single value by row and column label.

## Arithmetic and data alignment

One of the most important pandas features is the behavior of arithmetic between objects with different indexes. When adding together objects, if any index pairs are not the same, the respective index in the result will be the union of the index pairs. Let's look at a simple example:

```
In [126]: s1 = Series([7.3, -2.5, 3.4, 1.5], index=['a', 'c', 'd', 'e'])
```

```
In [127]: s2 = Series([-2.1, 3.6, -1.5, 4, 3.1], index=['a', 'c', 'e', 'f', 'g'])
```

In [128]: s1	In [129]: s2
Out[128]:	Out[129]:
a 7.3	a -2.1
c -2.5	c 3.6
d 3.4	e -1.5

```
e    1.5          f    4.0  
g    3.1
```

Adding these together yields:

```
In [130]: s1 + s2  
Out[130]:  
a    5.2  
c    1.1  
d    NaN  
e    0.0  
f    NaN  
g    NaN
```

The internal data alignment introduces NA values in the indices that don't overlap. Missing values propagate in arithmetic computations.

In the case of DataFrame, alignment is performed on both the rows and the columns:

```
In [131]: df1 = DataFrame(np.arange(9.).reshape((3, 3)), columns=list('bcd'),  
.....:                  index=['Ohio', 'Texas', 'Colorado'])  
  
In [132]: df2 = DataFrame(np.arange(12.).reshape((4, 3)), columns=list('bde'),  
.....:                  index=['Utah', 'Ohio', 'Texas', 'Oregon'])  
  
In [133]: df1  
Out[133]:  
      b   c   d  
Ohio    0   1   2  
Texas   3   4   5  
Colorado 6   7   8  
  
In [134]: df2  
Out[134]:  
      b   d   e  
Utah    0   1   2  
Ohio    3   4   5  
Texas   6   7   8  
Oregon  9  10  11
```

Adding these together returns a DataFrame whose index and columns are the unions of the ones in each DataFrame:

```
In [135]: df1 + df2  
Out[135]:  
      b   c   d   e  
Colorado  NaN  NaN  NaN  NaN  
Ohio     3   NaN  6   NaN  
Oregon   NaN  NaN  NaN  NaN  
Texas    9   NaN  12  NaN  
Utah    NaN  NaN  NaN  NaN
```

## Arithmetic methods with fill values

In arithmetic operations between differently-indexed objects, you might want to fill with a special value, like 0, when an axis label is found in one object but not the other:

```
In [136]: df1 = DataFrame(np.arange(12.).reshape((3, 4)), columns=list('abcd'))  
  
In [137]: df2 = DataFrame(np.arange(20.).reshape((4, 5)), columns=list('abcde'))  
  
In [138]: df1  
Out[138]:  
      a   b   c   d  
.....:  
.....:
```

```
In [139]: df2  
Out[139]:  
      a   b   c   d   e  
.....:  
.....:
```

```

0 0 1 2 3      0 0 1 2 3 4
1 4 5 6 7      1 5 6 7 8 9
2 8 9 10 11    2 10 11 12 13 14
                           3 15 16 17 18 19

```

Adding these together results in NA values in the locations that don't overlap:

In [140]: `df1 + df2`

Out[140]:

	a	b	c	d	e
0	0	2	4	6	NaN
1	9	11	13	15	NaN
2	18	20	22	24	NaN
3	NaN	NaN	NaN	NaN	NaN

Using the `add` method on `df1`, I pass `df2` and an argument to `fill_value`:

In [141]: `df1.add(df2, fill_value=0)`

Out[141]:

	a	b	c	d	e
0	0	2	4	6	4
1	9	11	13	15	9
2	18	20	22	24	14
3	15	16	17	18	19

Relatedly, when reindexing a Series or DataFrame, you can also specify a different fill value:

In [142]: `df1.reindex(columns=df2.columns, fill_value=0)`

Out[142]:

	a	b	c	d	e
0	0	1	2	3	0
1	4	5	6	7	0
2	8	9	10	11	0

*Table 5-7. Flexible arithmetic methods*

Method	Description
<code>add</code>	Method for addition (+)
<code>sub</code>	Method for subtraction (-)
<code>div</code>	Method for division (/)
<code>mul</code>	Method for multiplication (*)

## Operations between DataFrame and Series

As with NumPy arrays, arithmetic between DataFrame and Series is well-defined. First, as a motivating example, consider the difference between a 2D array and one of its rows:

In [143]: `arr = np.arange(12.).reshape((3, 4))`

In [144]: `arr`

Out[144]:

```

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

```

```

[ 8.,  9., 10., 11.]])

In [145]: arr[0]
Out[145]: array([ 0.,  1.,  2.,  3.])

In [146]: arr - arr[0]
Out[146]:
array([[ 0.,  0.,  0.,  0.],
       [ 4.,  4.,  4.,  4.],
       [ 8.,  8.,  8.,  8.]])

```

This is referred to as *broadcasting* and is explained in more detail in [Chapter 12](#). Operations between a DataFrame and a Series are similar:

```

In [147]: frame = DataFrame(np.arange(12.).reshape((4, 3)), columns=list('bde'),
                           index=['Utah', 'Ohio', 'Texas', 'Oregon'])

In [148]: series = frame.ix[0]

In [149]: frame           In [150]: series
Out[149]:              Out[150]:
      b   d   e           b   0
Utah  0   1   2           d   1
Ohio   3   4   5           e   2
Texas  6   7   8           Name: Utah
Oregon 9  10  11

```

By default, arithmetic between DataFrame and Series matches the index of the Series on the DataFrame's columns, broadcasting down the rows:

```

In [151]: frame - series
Out[151]:
      b   d   e
Utah  0   0   0
Ohio   3   3   3
Texas  6   6   6
Oregon 9  9  9

```

If an index value is not found in either the DataFrame's columns or the Series's index, the objects will be reindexed to form the union:

```

In [152]: series2 = Series(range(3), index=['b', 'e', 'f'])

In [153]: frame + series2
Out[153]:
      b   d   e   f
Utah  0  NaN  3  NaN
Ohio   3  NaN  6  NaN
Texas  6  NaN  9  NaN
Oregon 9  NaN 12  NaN

```

If you want to instead broadcast over the columns, matching on the rows, you have to use one of the arithmetic methods. For example:

```

In [154]: series3 = frame['d']

In [155]: frame      In [156]: series3

```

```

Out[155]:          Out[156]:
      b  d  e    Utah     1
Utah  0  1  2    Ohio     4
Ohio  3  4  5   Texas     7
Texas 6  7  8  Oregon    10
Oregon 9 10 11  Name: d

In [157]: frame.sub(series3, axis=0)
Out[157]:
      b  d  e
Utah -1  0  1
Ohio -1  0  1
Texas -1  0  1
Oregon -1  0  1

```

The axis number that you pass is the *axis to match on*. In this case we mean to match on the DataFrame's row index and broadcast across.

## Function application and mapping

NumPy ufuncs (element-wise array methods) work fine with pandas objects:

```

In [158]: frame = DataFrame(np.random.randn(4, 3), columns=list('bde'),
.....:                  index=['Utah', 'Ohio', 'Texas', 'Oregon'])

In [159]: frame
Out[159]:
      b         d         e
Utah -0.204708  0.478943 -0.519439
Ohio -0.555730  1.965781  1.393406
Texas 0.092908  0.281746  0.769023
Oregon 1.246435  1.007189 -1.296221

In [160]: np.abs(frame)
Out[160]:
      b         d         e
Utah  0.204708  0.478943  0.519439
Ohio  0.555730  1.965781  1.393406
Texas 0.092908  0.281746  0.769023
Oregon 1.246435  1.007189  1.296221

```

Another frequent operation is applying a function on 1D arrays to each column or row. DataFrame's `apply` method does exactly this:

```

In [161]: f = lambda x: x.max() - x.min()

In [162]: frame.apply(f)
Out[162]:
      b        d        e
b  1.802165
d  1.684034
e  2.689627

In [163]: frame.apply(f, axis=1)
Out[163]:
      Utah    Ohio    Texas    Oregon
0.998382 2.521511 0.676115 2.542656

```

Many of the most common array statistics (like `sum` and `mean`) are DataFrame methods, so using `apply` is not necessary.

The function passed to `apply` need not return a scalar value, it can also return a Series with multiple values:

```

In [164]: def f(x):
.....:     return Series([x.min(), x.max()], index=['min', 'max'])

In [165]: frame.apply(f)

```

```
Out[165]:  
      b      d      e  
min -0.555730  0.281746 -1.296221  
max  1.246435  1.965781  1.393406
```

Element-wise Python functions can be used, too. Suppose you wanted to compute a formatted string from each floating point value in `frame`. You can do this with `applymap`:

```
In [166]: format = lambda x: '%.2f' % x
```

```
In [167]: frame.applymap(format)  
Out[167]:
```

	b	d	e
Utah	-0.20	0.48	-0.52
Ohio	-0.56	1.97	1.39
Texas	0.09	0.28	0.77
Oregon	1.25	1.01	-1.30

The reason for the name `applymap` is that Series has a `map` method for applying an element-wise function:

```
In [168]: frame['e'].map(format)  
Out[168]:  
Utah      -0.52  
Ohio       1.39  
Texas      0.77  
Oregon     -1.30  
Name: e
```

## Sorting and ranking

Sorting a data set by some criterion is another important built-in operation. To sort lexicographically by row or column index, use the `sort_index` method, which returns a new, sorted object:

```
In [169]: obj = Series(range(4), index=['d', 'a', 'b', 'c'])
```

```
In [170]: obj.sort_index()  
Out[170]:  
a    1  
b    2  
c    3  
d    0
```

With a DataFrame, you can sort by index on either axis:

```
In [171]: frame = DataFrame(np.arange(8).reshape((2, 4)), index=['three', 'one'],  
.....:           columns=['d', 'a', 'b', 'c'])
```

```
In [172]: frame.sort_index()  
Out[172]:  
      d  a  b  c  
one   4  5  6  7  
three  0  1  2  3
```

```
In [173]: frame.sort_index(axis=1)  
Out[173]:  
      a  b  c  d  
three  1  2  3  0  
one    5  6  7  4
```

The data is sorted in ascending order by default, but can be sorted in descending order, too:

```
In [174]: frame.sort_index(axis=1, ascending=False)
Out[174]:
      d  c  b  a
three  0  3  2  1
one    4  7  6  5
```

To sort a Series by its values, use its `order` method:

```
In [175]: obj = Series([4, 7, -3, 2])
In [176]: obj.order()
Out[176]:
2    -3
3     2
0     4
1     7
```

Any missing values are sorted to the end of the Series by default:

```
In [177]: obj = Series([4, np.nan, 7, np.nan, -3, 2])
In [178]: obj.order()
Out[178]:
4    -3
5     2
0     4
2     7
1    NaN
3    NaN
```

On DataFrame, you may want to sort by the values in one or more columns. To do so, pass one or more column names to the `by` option:

```
In [179]: frame = DataFrame({'b': [4, 7, -3, 2], 'a': [0, 1, 0, 1]})
In [180]: frame           In [181]: frame.sort_index(by='b')
Out[180]:               Out[181]:
      a   b                 a   b
0   0   4                 2   0  -3
1   1   7                 3   1   2
2   0  -3                 0   0   4
3   1   2                 1   1   7
```

To sort by multiple columns, pass a list of names:

```
In [182]: frame.sort_index(by=['a', 'b'])
Out[182]:
      a   b
2   0  -3
0   0   4
3   1   2
1   1   7
```

*Ranking* is closely related to sorting, assigning ranks from one through the number of valid data points in an array. It is similar to the indirect sort indices produced by `numpy.argsort`, except that ties are broken according to a rule. The `rank` methods for Series and DataFrame are the place to look; by default `rank` breaks ties by assigning each group the mean rank:

```
In [183]: obj = Series([7, -5, 7, 4, 2, 0, 4])
```

```
In [184]: obj.rank()
```

```
Out[184]:
```

0	6.5
1	1.0
2	6.5
3	4.5
4	3.0
5	2.0
6	4.5

Ranks can also be assigned according to the order they're observed in the data:

```
In [185]: obj.rank(method='first')
```

```
Out[185]:
```

0	6
1	1
2	7
3	4
4	3
5	2
6	5

Naturally, you can rank in descending order, too:

```
In [186]: obj.rank(ascending=False, method='max')
```

```
Out[186]:
```

0	2
1	7
2	2
3	4
4	5
5	6
6	4

See [Table 5-8](#) for a list of tie-breaking methods available. DataFrame can compute ranks over the rows or the columns:

```
In [187]: frame = DataFrame({'b': [4.3, 7, -3, 2], 'a': [0, 1, 0, 1],  
.....: 'c': [-2, 5, 8, -2.5]})
```

```
In [188]: frame
```

```
Out[188]:
```

	a	b	c
0	0	4.3	-2.0
1	1	7.0	5.0
2	0	-3.0	8.0
3	1	2.0	-2.5

```
In [189]: frame.rank(axis=1)
```

```
Out[189]:
```

	a	b	c
0	2	3	1
1	1	3	2
2	2	1	3
3	2	3	1

Table 5-8. Tie-breaking methods with rank

Method	Description
'average'	Default: assign the average rank to each entry in the equal group.
'min'	Use the minimum rank for the whole group.
'max'	Use the maximum rank for the whole group.
'first'	Assign ranks in the order the values appear in the data.

## Axis indexes with duplicate values

Up until now all of the examples I've showed you have had unique axis labels (index values). While many pandas functions (like `reindex`) require that the labels be unique, it's not mandatory. Let's consider a small Series with duplicate indices:

```
In [190]: obj = Series(range(5), index=['a', 'a', 'b', 'b', 'c'])
```

```
In [191]: obj
```

```
Out[191]:
```

```
a    0  
a    1  
b    2  
b    3  
c    4
```

The index's `is_unique` property can tell you whether its values are unique or not:

```
In [192]: obj.index.is_unique
```

```
Out[192]: False
```

Data selection is one of the main things that behaves differently with duplicates. Indexing a value with multiple entries returns a Series while single entries return a scalar value:

```
In [193]: obj['a']      In [194]: obj['c']  
Out[193]:                  Out[194]: 4  
a    0  
a    1
```

The same logic extends to indexing rows in a DataFrame:

```
In [195]: df = DataFrame(np.random.randn(4, 3), index=['a', 'a', 'b', 'b'])
```

```
In [196]: df
```

```
Out[196]:
```

```
          0         1         2  
a  0.274992  0.228913  1.352917  
a  0.886429 -2.001637 -0.371843  
b  1.669025 -0.438570 -0.539741  
b  0.476985  3.248944 -1.021228
```

```
In [197]: df.ix['b']
```

```
Out[197]:
```

```
          0         1         2
```

```
b  1.669025 -0.438570 -0.539741  
b  0.476985  3.248944 -1.021228
```

## Summarizing and Computing Descriptive Statistics

pandas objects are equipped with a set of common mathematical and statistical methods. Most of these fall into the category of *reductions* or *summary statistics*, methods that extract a single value (like the sum or mean) from a Series or a Series of values from the rows or columns of a DataFrame. Compared with the equivalent methods of vanilla NumPy arrays, they are all built from the ground up to exclude missing data. Consider a small DataFrame:

```
In [198]: 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'])  
  
In [199]: df  
Out[199]:  
      one   two  
a  1.40  NaN  
b  7.10 -4.5  
c  NaN   NaN  
d  0.75 -1.3
```

Calling DataFrame's `sum` method returns a Series containing column sums:

```
In [200]: df.sum()  
Out[200]:  
one    9.25  
two   -5.80
```

Passing `axis=1` sums over the rows instead:

```
In [201]: df.sum(axis=1)  
Out[201]:  
a    1.40  
b    2.60  
c    NaN  
d   -0.55
```

NA values are excluded unless the entire slice (row or column in this case) is NA. This can be disabled using the `skipna` option:

```
In [202]: df.mean(axis=1, skipna=False)  
Out[202]:  
a    NaN  
b    1.300  
c    NaN  
d   -0.275
```

See [Table 5-9](#) for a list of common options for each reduction method options.

Table 5-9. Options for reduction methods

Method	Description
axis	Axis to reduce over. 0 for DataFrame's rows and 1 for columns.
skipna	Exclude missing values, True by default.
level	Reduce grouped by level if the axis is hierarchically-indexed (MultiIndex).

Some methods, like `idxmin` and `idxmax`, return indirect statistics like the index value where the minimum or maximum values are attained:

```
In [203]: df.idxmax()  
Out[203]:  
one    b  
two    d
```

Other methods are *accumulations*:

```
In [204]: df.cumsum()  
Out[204]:  
      one   two  
a  1.40  NaN  
b  8.50 -4.5  
c  NaN   NaN  
d  9.25 -5.8
```

Another type of method is neither a reduction nor an accumulation. `describe` is one such example, producing multiple summary statistics in one shot:

```
In [205]: df.describe()  
Out[205]:  
      one      two  
count  3.000000  2.000000  
mean   3.083333 -2.900000  
std    3.493685  2.262742  
min    0.750000 -4.500000  
25%   1.075000 -3.700000  
50%   1.400000 -2.900000  
75%   4.250000 -2.100000  
max    7.100000 -1.300000
```

On non-numeric data, `describe` produces alternate summary statistics:

```
In [206]: obj = Series(['a', 'a', 'b', 'c'] * 4)  
  
In [207]: obj.describe()  
Out[207]:  
count      16  
unique      3  
top        a  
freq       8
```

See [Table 5-10](#) for a full list of summary statistics and related methods.

Table 5-10. Descriptive and summary statistics

Method	Description
count	Number of non-NA values
describe	Compute set of summary statistics for Series or each DataFrame column
min, max	Compute minimum and maximum values
argmin, argmax	Compute index locations (integers) at which minimum or maximum value obtained, respectively
idxmin, idxmax	Compute index values at which minimum or maximum value obtained, respectively
quantile	Compute sample quantile ranging from 0 to 1
sum	Sum of values
mean	Mean of values
median	Arithmetic median (50% quantile) of values
mad	Mean absolute deviation from mean value
var	Sample variance of values
std	Sample standard deviation of values
skew	Sample skewness (3rd moment) of values
kurt	Sample kurtosis (4th moment) of values
cumsum	Cumulative sum of values
cummin, cummax	Cumulative minimum or maximum of values, respectively
cumprod	Cumulative product of values
diff	Compute 1st arithmetic difference (useful for time series)
pct_change	Compute percent changes

## Correlation and Covariance

Some summary statistics, like correlation and covariance, are computed from pairs of arguments. Let's consider some DataFrames of stock prices and volumes obtained from Yahoo! Finance:

```
import pandas.io.data as web

all_data = {}
for ticker in ['AAPL', 'IBM', 'MSFT', 'GOOG']:
    all_data[ticker] = web.get_data_yahoo(ticker, '1/1/2000', '1/1/2010')

price = DataFrame({tic: data['Adj Close']
                  for tic, data in all_data.iteritems()})
volume = DataFrame({tic: data['Volume']
                   for tic, data in all_data.iteritems()})
```

I now compute percent changes of the prices:

```
In [209]: returns = price.pct_change()
In [210]: returns.tail()
```

```
Out[210]:
```

	AAPL	GOOG	IBM	MSFT
Date				
2009-12-24	0.034339	0.011117	0.004420	0.002747
2009-12-28	0.012294	0.007098	0.013282	0.005479
2009-12-29	-0.011861	-0.005571	-0.003474	0.006812
2009-12-30	0.012147	0.005376	0.005468	-0.013532
2009-12-31	-0.004300	-0.004416	-0.012609	-0.015432

The `corr` method of Series computes the correlation of the overlapping, non-NA, aligned-by-index values in two Series. Relatedly, `cov` computes the covariance:

```
In [211]: returns.MSFT.corr(returns.IBM)
```

```
Out[211]: 0.49609291822168838
```

```
In [212]: returns.MSFT.cov(returns.IBM)
```

```
Out[212]: 0.00021600332437329015
```

DataFrame's `corr` and `cov` methods, on the other hand, return a full correlation or covariance matrix as a DataFrame, respectively:

```
In [213]: returns.corr()
```

```
Out[213]:
```

	AAPL	GOOG	IBM	MSFT
AAPL	1.000000	0.470660	0.410648	0.424550
GOOG	0.470660	1.000000	0.390692	0.443334
IBM	0.410648	0.390692	1.000000	0.496093
MSFT	0.424550	0.443334	0.496093	1.000000

```
In [214]: returns.cov()
```

```
Out[214]:
```

	AAPL	GOOG	IBM	MSFT
AAPL	0.001028	0.000303	0.000252	0.000309
GOOG	0.000303	0.000580	0.000142	0.000205
IBM	0.000252	0.000142	0.000367	0.000216
MSFT	0.000309	0.000205	0.000216	0.000516

Using DataFrame's `corrwith` method, you can compute pairwise correlations between a DataFrame's columns or rows with another Series or DataFrame. Passing a Series returns a Series with the correlation value computed for each column:

```
In [215]: returns.corrwith(returns.IBM)
```

```
Out[215]:
```

AAPL	0.410648
GOOG	0.390692
IBM	1.000000
MSFT	0.496093

Passing a DataFrame computes the correlations of matching column names. Here I compute correlations of percent changes with volume:

```
In [216]: returns.corrwith(volume)
```

```
Out[216]:
```

AAPL	-0.057461
GOOG	0.062644

```
IBM    -0.007900
MSFT   -0.014175
```

Passing `axis=1` does things row-wise instead. In all cases, the data points are aligned by label before computing the correlation.

## Unique Values, Value Counts, and Membership

Another class of related methods extracts information about the values contained in a one-dimensional Series. To illustrate these, consider this example:

```
In [217]: obj = Series(['c', 'a', 'd', 'a', 'a', 'b', 'b', 'c', 'c'])
```

The first function is `unique`, which gives you an array of the unique values in a Series:

```
In [218]: uniques = obj.unique()
```

```
In [219]: uniques
```

```
Out[219]: array(['c', 'a', 'd', 'b'], dtype=object)
```

The unique values are not necessarily returned in sorted order, but could be sorted after the fact if needed (`uniques.sort()`). Relatedly, `value_counts` computes a Series containing value frequencies:

```
In [220]: obj.value_counts()
```

```
Out[220]:
```

```
c    3
a    3
b    2
d    1
```

The Series is sorted by value in descending order as a convenience. `value_counts` is also available as a top-level pandas method that can be used with any array or sequence:

```
In [221]: pd.value_counts(obj.values, sort=False)
```

```
Out[221]:
```

```
a    3
b    2
c    3
d    1
```

Lastly, `isin` is responsible for vectorized set membership and can be very useful in filtering a data set down to a subset of values in a Series or column in a DataFrame:

```
In [222]: mask = obj.isin(['b', 'c'])
```

```
In [223]: mask
```

```
Out[223]:
```

```
0    True
1   False
2   False
3   False
4   False
5    True
6    True
```

```
In [224]: obj[mask]
```

```
Out[224]:
```

```
0    c
5    b
6    b
7    c
8    c
```

```
7    True
8    True
```

See [Table 5-11](#) for a reference on these methods.

*Table 5-11. Unique, value counts, and binning methods*

Method	Description
isin	Compute boolean array indicating whether each Series value is contained in the passed sequence of values.
unique	Compute array of unique values in a Series, returned in the order observed.
value_counts	Return a Series containing unique values as its index and frequencies as its values, ordered count in descending order.

In some cases, you may want to compute a histogram on multiple related columns in a DataFrame. Here's an example:

```
In [225]: data = DataFrame({'Qu1': [1, 3, 4, 3, 4],
.....:                 'Qu2': [2, 3, 1, 2, 3],
.....:                 'Qu3': [1, 5, 2, 4, 4]})

In [226]: data
Out[226]:
   Qu1  Qu2  Qu3
0     1     2     1
1     3     3     5
2     4     1     2
3     3     2     4
4     4     3     4
```

Passing `pandas.value_counts` to this DataFrame's `apply` function gives:

```
In [227]: result = data.apply(pd.value_counts).fillna(0)

In [228]: result
Out[228]:
   Qu1  Qu2  Qu3
1     1     1     1
2     0     2     1
3     2     2     0
4     2     0     2
5     0     0     1
```

## Handling Missing Data

Missing data is common in most data analysis applications. One of the goals in designing pandas was to make working with missing data as painless as possible. For example, all of the descriptive statistics on pandas objects exclude missing data as you've seen earlier in the chapter.

pandas uses the floating point value `NaN` (Not a Number) to represent missing data in both floating as well as in non-floating point arrays. It is just used as a *sentinel* that can be easily detected:

```
In [229]: string_data = Series(['aardvark', 'artichoke', np.nan, 'avocado'])

In [230]: string_data
Out[230]:
0    aardvark
1    artichoke
2      NaN
3    avocado

In [231]: string_data.isnull()
Out[231]:
0    False
1    False
2     True
3    False
```

The built-in Python `None` value is also treated as NA in object arrays:

```
In [232]: string_data[0] = None

In [233]: string_data.isnull()
Out[233]:
0     True
1    False
2     True
3    False
```

I do not claim that pandas's NA representation is optimal, but it is simple and reasonably consistent. It's the best solution, with good all-around performance characteristics and a simple API, that I could concoct in the absence of a true NA data type or bit pattern in NumPy's data types. Ongoing development work in NumPy may change this in the future.

Table 5-12. NA handling methods

Argument	Description
<code>dropna</code>	Filter axis labels based on whether values for each label have missing data, with varying thresholds for how much missing data to tolerate.
<code>fillna</code>	Fill in missing data with some value or using an interpolation method such as ' <code>ffill</code> ' or ' <code>bfill</code> '.
<code>isnull</code>	Return like-type object containing boolean values indicating which values are missing / NA.
<code>notnull</code>	Negation of <code>isnull</code> .

## Filtering Out Missing Data

You have a number of options for filtering out missing data. While doing it by hand is always an option, `dropna` can be very helpful. On a Series, it returns the Series with only the non-null data and index values:

```
In [234]: from numpy import nan as NA

In [235]: data = Series([1, NA, 3.5, NA, 7])

In [236]: data.dropna()
Out[236]:
```

```
0    1.0  
2    3.5  
4    7.0
```

Naturally, you could have computed this yourself by boolean indexing:

```
In [237]: data[data.notnull()]  
Out[237]:  
0    1.0  
2    3.5  
4    7.0
```

With DataFrame objects, these are a bit more complex. You may want to drop rows or columns which are all NA or just those containing any NAs. `dropna` by default drops any row containing a missing value:

```
In [238]: data = DataFrame([[1., 6.5, 3.], [1., NA, NA],  
.....:                 [NA, NA, NA], [NA, 6.5, 3.]])
```

```
In [239]: cleaned = data.dropna()
```

```
In [240]: data           In [241]: cleaned  
Out[240]:          Out[241]:  
0   1   2           0   1   2  
0   1   6.5   3     0   1   6.5   3  
1   1   NaN  NaN  
2  NaN  NaN  NaN  
3  NaN   6.5   3
```

Passing `how='all'` will only drop rows that are all NA:

```
In [242]: data.dropna(how='all')  
Out[242]:  
0   1   2  
0   1   6.5   3  
1   1   NaN  NaN  
3  NaN   6.5   3
```

Dropping columns in the same way is only a matter of passing `axis=1`:

```
In [243]: data[4] = NA
```

```
In [244]: data           In [245]: data.dropna(axis=1, how='all')  
Out[244]:          Out[245]:  
0   1   2   4           0   1   2  
0   1   6.5   3  NaN     0   1   6.5   3  
1   1   NaN  NaN  NaN     1   1   NaN  NaN  
2  NaN  NaN  NaN  NaN     2  NaN  NaN  NaN  
3  NaN   6.5   3  NaN     3  NaN   6.5   3
```

A related way to filter out DataFrame rows tends to concern time series data. Suppose you want to keep only rows containing a certain number of observations. You can indicate this with the `thresh` argument:

```
In [246]: df = DataFrame(np.random.randn(7, 3))
```

```
In [247]: df.ix[:4, 1] = NA; df.ix[:2, 2] = NA
```

```
In [248]: df  
Out[248]:
```

	0	1	2
0	-0.577087	NaN	NaN
1	0.523772	NaN	NaN
2	-0.713544	NaN	NaN
3	-1.860761	NaN	0.560145
4	-1.265934	NaN	-1.063512
5	0.332883	-2.359419	-0.199543
6	-1.541996	-0.970736	-1.307030

```
In [249]: df.dropna(thresh=3)  
Out[249]:
```

	0	1	2
5	0.332883	-2.359419	-0.199543
6	-1.541996	-0.970736	-1.307030

## Filling in Missing Data

Rather than filtering out missing data (and potentially discarding other data along with it), you may want to fill in the “holes” in any number of ways. For most purposes, the `fillna` method is the workhorse function to use. Calling `fillna` with a constant replaces missing values with that value:

```
In [250]: df.fillna(0)  
Out[250]:
```

	0	1	2
0	-0.577087	0.000000	0.000000
1	0.523772	0.000000	0.000000
2	-0.713544	0.000000	0.000000
3	-1.860761	0.000000	0.560145
4	-1.265934	0.000000	-1.063512
5	0.332883	-2.359419	-0.199543
6	-1.541996	-0.970736	-1.307030

Calling `fillna` with a dict you can use a different fill value for each column:

```
In [251]: df.fillna({1: 0.5, 3: -1})  
Out[251]:
```

	0	1	2
0	-0.577087	0.500000	NaN
1	0.523772	0.500000	NaN
2	-0.713544	0.500000	NaN
3	-1.860761	0.500000	0.560145
4	-1.265934	0.500000	-1.063512
5	0.332883	-2.359419	-0.199543
6	-1.541996	-0.970736	-1.307030

`fillna` returns a new object, but you can modify the existing object in place:

```
# always returns a reference to the filled object  
In [252]: _ = df.fillna(0, inplace=True)
```

```
In [253]: df  
Out[253]:
```

	0	1	2
0	-0.577087	0.000000	0.000000
1	0.523772	0.000000	0.000000
2	-0.713544	0.000000	0.000000
3	-1.860761	0.000000	0.560145

```
4 -1.265934  0.000000 -1.063512
5  0.332883 -2.359419 -0.199543
6 -1.541996 -0.970736 -1.307030
```

The same interpolation methods available for reindexing can be used with `fillna`:

```
In [254]: df = DataFrame(np.random.randn(6, 3))
```

```
In [255]: df.ix[2:, 1] = NA; df.ix[4:, 2] = NA
```

```
In [256]: df
```

```
Out[256]:
```

```
      0         1         2
0  0.286350  0.377984 -0.753887
1  0.331286  1.349742  0.069877
2  0.246674    NaN  1.004812
3  1.327195    NaN -1.549106
4  0.022185    NaN     NaN
5  0.862580    NaN     NaN
```

```
In [257]: df.fillna(method='ffill')
```

```
Out[257]:
```

```
      0         1         2
0  0.286350  0.377984 -0.753887
1  0.331286  1.349742  0.069877
2  0.246674  1.349742  1.004812
3  1.327195  1.349742 -1.549106
4  0.022185  1.349742 -1.549106
5  0.862580  1.349742 -1.549106
```

```
In [258]: df.fillna(method='ffill', limit=2)
```

```
Out[258]:
```

```
      0         1         2
0  0.286350  0.377984 -0.753887
1  0.331286  1.349742  0.069877
2  0.246674  1.349742  1.004812
3  1.327195  1.349742 -1.549106
4  0.022185    NaN -1.549106
5  0.862580    NaN -1.549106
```

With `fillna` you can do lots of other things with a little creativity. For example, you might pass the mean or median value of a Series:

```
In [259]: data = Series([1., NA, 3.5, NA, 7])
```

```
In [260]: data.fillna(data.mean())
```

```
Out[260]:
```

```
0    1.000000
1    3.833333
2    3.500000
3    3.833333
4    7.000000
```

See [Table 5-13](#) for a reference on `fillna`.

*Table 5-13. `fillna` function arguments*

Argument	Description
<code>value</code>	Scalar value or dict-like object to use to fill missing values
<code>method</code>	Interpolation, by default ' <code>ffill</code> ' if function called with no other arguments
<code>axis</code>	Axis to fill on, default <code>axis=0</code>
<code>inplace</code>	Modify the calling object without producing a copy
<code>limit</code>	For forward and backward filling, maximum number of consecutive periods to fill

# Hierarchical Indexing

*Hierarchical indexing* is an important feature of pandas enabling you to have multiple (two or more) index *levels* on an axis. Somewhat abstractly, it provides a way for you to work with higher dimensional data in a lower dimensional form. Let's start with a simple example; create a Series with a list of lists or arrays as the index:

```
In [261]: data = Series(np.random.randn(10),
.....:                 index=[['a', 'a', 'a', 'b', 'b', 'b', 'c', 'c', 'd', 'd'],
.....:                     [1, 2, 3, 1, 2, 3, 1, 2, 2, 3]])
```

```
In [262]: data
Out[262]:
a    0.670216
     0.852965
     -0.955869
b    -0.023493
     -2.304234
     -0.652469
c    -1.218302
     -1.332610
d    1.074623
     0.723642
```

What you're seeing is a prettified view of a Series with a `MultiIndex` as its index. The “gaps” in the index display mean “use the label directly above”:

```
In [263]: data.index
Out[263]:
MultiIndex
[('a', 1) ('a', 2) ('a', 3) ('b', 1) ('b', 2) ('b', 3) ('c', 1)
 ('c', 2) ('d', 2) ('d', 3)]
```

With a hierarchically-indexed object, so-called *partial* indexing is possible, enabling you to concisely select subsets of the data:

```
In [264]: data['b']
Out[264]:
1    -0.023493
2    -2.304234
3    -0.652469
```

```
In [265]: data['b':'c']
Out[265]:
b    -0.023493
     -2.304234
     -0.652469
c    -1.218302
     -1.332610
```

```
In [266]: data.ix[['b', 'd']]
Out[266]:
b    -0.023493
     -2.304234
     -0.652469
d    1.074623
     0.723642
```

Selection is even possible in some cases from an “inner” level:

```
In [267]: data[:, 2]
Out[267]:
a    0.852965
```

```
b    -2.304234  
c    -1.332610  
d     1.074623
```

Hierarchical indexing plays a critical role in reshaping data and group-based operations like forming a pivot table. For example, this data could be rearranged into a DataFrame using its `unstack` method:

```
In [268]: data.unstack()  
Out[268]:  
      1         2         3  
a  0.670216  0.852965 -0.955869  
b -0.023493 -2.304234 -0.652469  
c -1.218302 -1.332610      NaN  
d      NaN  1.074623  0.723642
```

The inverse operation of `unstack` is `stack`:

```
In [269]: data.unstack().stack()  
Out[269]:  
a  1    0.670216  
   2    0.852965  
   3   -0.955869  
b  1   -0.023493  
   2   -2.304234  
   3   -0.652469  
c  1   -1.218302  
   2   -1.332610  
   3    1.074623  
      3    0.723642
```

`stack` and `unstack` will be explored in more detail in [Chapter 7](#).

With a DataFrame, either axis can have a hierarchical index:

```
In [270]: frame = DataFrame(np.arange(12).reshape((4, 3)),  
.....:                      index=[['a', 'a', 'b', 'b'], [1, 2, 1, 2]],  
.....:                      columns=[['Ohio', 'Ohio', 'Colorado'],  
.....:                                ['Green', 'Red', 'Green']])  
  
In [271]: frame  
Out[271]:  
          Ohio        Colorado  
          Green   Red   Green  
a 1      0     1     2  
   2      3     4     5  
b 1      6     7     8  
   2      9    10    11
```

The hierarchical levels can have names (as strings or any Python objects). If so, these will show up in the console output (don't confuse the index names with the axis labels!):

```
In [272]: frame.index.names = ['key1', 'key2']  
  
In [273]: frame.columns.names = ['state', 'color']  
  
In [274]: frame
```

```

Out[274]:
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

```

With partial column indexing you can similarly select groups of columns:

```

In [275]: frame['Ohio']
Out[275]:
color     Green   Red
key1 key2
a      1        0        1
      2        3        4
b      1        6        7
      2        9       10

```

A MultiIndex can be created by itself and then reused; the columns in the above DataFrame with level names could be created like this:

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

## Reordering and Sorting Levels

At times you will need to rearrange the order of the levels on an axis or sort the data by the values in one specific level. The `swaplevel` takes two level numbers or names and returns a new object with the levels interchanged (but the data is otherwise unaltered):

```

In [276]: frame.swaplevel('key1', 'key2')
Out[276]:
state      Ohio      Colorado
color     Green    Red     Green
key2 key1
1      a        0        1        2
2      a        3        4        5
1      b        6        7        8
2      b        9       10       11

```

`sortlevel`, on the other hand, sorts the data (stably) using only the values in a single level. When swapping levels, it's not uncommon to also use `sortlevel` so that the result is lexicographically sorted:

```

In [277]: frame.sortlevel(1)
Out[277]:
state      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 [278]: frame.swaplevel(0, 1).sortlevel(0)
Out[278]:
state      Ohio      Colorado
color     Green    Red     Green
key2 key1
1      a        0        1        2
b      a        3        4        5
2      a        6        7        8
b      b        9       10       11

```



Data selection performance is much better on hierarchically indexed objects if the index is lexicographically sorted starting with the outermost level, that is, the result of calling `sortlevel(0)` or `sort_index()`.

## Summary Statistics by Level

Many descriptive and summary statistics on DataFrame and Series have a `level` option in which you can specify the level you want to sum by on a particular axis. Consider the above DataFrame; we can sum by level on either the rows or columns like so:

```
In [279]: frame.sum(level='key2')
Out[279]:
state    Ohio      Colorado
color   Green    Red      Green
key2
1          6     8       10
2         12    14       16

In [280]: frame.sum(level='color', axis=1)
Out[280]:
color   Green  Red
key1 key2
a      1     2     1
      2     8     4
b      1    14     7
      2    20    10
```

Under the hood, this utilizes pandas's `groupby` machinery which will be discussed in more detail later in the book.

## Using a DataFrame's Columns

It's not unusual to want to use one or more columns from a DataFrame as the row index; alternatively, you may wish to move the row index into the DataFrame's columns. Here's an example DataFrame:

```
In [281]: frame = DataFrame({'a': range(7), 'b': range(7, 0, -1),
.....:                      'c': ['one', 'one', 'one', 'two', 'two', 'two'],
.....:                      'd': [0, 1, 2, 0, 1, 2, 3]})

In [282]: frame
Out[282]:
   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's `set_index` function will create a new DataFrame using one or more of its columns as the index:

```
In [283]: frame2 = frame.set_index(['c', 'd'])
```

```
In [284]: frame2
```

```
Out[284]:
```

		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

By default the columns are removed from the DataFrame, though you can leave them in:

```
In [285]: frame.set_index(['c', 'd'], drop=False)
```

```
Out[285]:
```

		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`, on the other hand, does the opposite of `set_index`; the hierarchical index levels are moved into the columns:

```
In [286]: frame2.reset_index()
```

```
Out[286]:
```

	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

## Other pandas Topics

Here are some additional topics that may be of use to you in your data travels.

### Integer Indexing

Working with pandas objects indexed by integers is something that often trips up new users due to some differences with indexing semantics on built-in Python data

structures like lists and tuples. For example, you would not expect the following code to generate an error:

```
ser = Series(np.arange(3.))
ser[-1]
```

In this case, pandas could “fall back” on integer indexing, but there’s not a safe and general way (that I know of) to do this without introducing subtle bugs. Here we have an index containing 0, 1, 2, but inferring what the user wants (label-based indexing or position-based) is difficult::

```
In [288]: ser
Out[288]:
0    0
1    1
2    2
```

On the other hand, with a non-integer index, there is no potential for ambiguity:

```
In [289]: ser2 = Series(np.arange(3.), index=['a', 'b', 'c'])

In [290]: ser2[-1]
Out[290]: 2.0
```

To keep things consistent, if you have an axis index containing indexers, data selection with integers will always be label-oriented. This includes slicing with `ix`, too:

```
In [291]: ser.ix[:1]
Out[291]:
0    0
1    1
```

In cases where you need reliable position-based indexing regardless of the index type, you can use the `iget_value` method from Series and `irow` and `icol` methods from DataFrame:

```
In [292]: ser3 = Series(range(3), index=[-5, 1, 3])

In [293]: ser3.iget_value(2)
Out[293]: 2

In [294]: frame = DataFrame(np.arange(6).reshape(3, 2), index=[2, 0, 1])

In [295]: frame.irow(0)
Out[295]:
0    0
1    1
Name: 2
```

## Panel Data

While not a major topic of this book, pandas has a Panel data structure, which you can think of as a three-dimensional analogue of DataFrame. Much of the development focus of pandas has been in tabular data manipulations as these are easier to reason about,

and hierarchical indexing makes using truly N-dimensional arrays unnecessary in a lot of cases.

To create a Panel, you can use a dict of DataFrame objects or a three-dimensional ndarray:

```
import pandas.io.data as web

pdata = pd.Panel(dict((stk, web.get_data_yahoo(stk, '1/1/2009', '6/1/2012'))
                      for stk in ['AAPL', 'GOOG', 'MSFT', 'DELL']))
```

Each item (the analogue of columns in a DataFrame) in the Panel is a DataFrame:

```
In [297]: pdata
Out[297]:
<class 'pandas.core.panel.Panel'>
Dimensions: 4 (items) x 861 (major) x 6 (minor)
Items: AAPL to MSFT
Major axis: 2009-01-02 00:00:00 to 2012-06-01 00:00:00
Minor axis: Open to Adj Close

In [298]: pdata = pdata.swapaxes('items', 'minor')

In [299]: pdata['Adj Close']
Out[299]:
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 861 entries, 2009-01-02 00:00:00 to 2012-06-01 00:00:00
Data columns:
AAPL    861 non-null values
DELL    861 non-null values
GOOG    861 non-null values
MSFT    861 non-null values
dtypes: float64(4)
```

ix-based label indexing generalizes to three dimensions, so we can select all data at a particular date or a range of dates like so:

```
In [300]: pdata.ix[:, '6/1/2012', :]
Out[300]:
      Open   High    Low  Close  Volume  Adj Close
AAPL  569.16  572.65  560.52  560.99  18606700    560.99
DELL   12.15   12.30   12.05   12.07  19396700    12.07
GOOG  571.79  572.65  568.35  570.98  3057900    570.98
MSFT   28.76   28.96   28.44   28.45  56634300    28.45

In [301]: pdata.ix['Adj Close', '5/22/2012':, :]
Out[301]:
          AAPL     DELL     GOOG     MSFT
Date
2012-05-22  556.97  15.08  600.80  29.76
2012-05-23  570.56  12.49  609.46  29.11
2012-05-24  565.32  12.45  603.66  29.07
2012-05-25  562.29  12.46  591.53  29.06
2012-05-29  572.27  12.66  594.34  29.56
2012-05-30  579.17  12.56  588.23  29.34
```

```
2012-05-31  577.73  12.33  580.86  29.19
2012-06-01  560.99  12.07  570.98  28.45
```

An alternate way to represent panel data, especially for fitting statistical models, is in “stacked” DataFrame form:

```
In [302]: stacked = pdata.ix[:, '5/30/2012':, :].to_frame()
```

```
In [303]: stacked
```

```
Out[303]:
```

		Open	High	Low	Close	Volume	Adj Close
major	minor						
2012-05-30	AAPL	569.20	579.99	566.56	579.17	18908200	579.17
	DELL	12.59	12.70	12.46	12.56	19787800	12.56
	GOOG	588.16	591.90	583.53	588.23	1906700	588.23
	MSFT	29.35	29.48	29.12	29.34	41585500	29.34
2012-05-31	AAPL	580.74	581.50	571.46	577.73	17559800	577.73
	DELL	12.53	12.54	12.33	12.33	19955500	12.33
	GOOG	588.72	590.00	579.00	580.86	2968300	580.86
	MSFT	29.30	29.42	28.94	29.19	39134000	29.19
2012-06-01	AAPL	569.16	572.65	560.52	560.99	18606700	560.99
	DELL	12.15	12.30	12.05	12.07	19396700	12.07
	GOOG	571.79	572.65	568.35	570.98	3057900	570.98
	MSFT	28.76	28.96	28.44	28.45	56634300	28.45

DataFrame has a related `to_panel` method, the inverse of `to_frame`:

```
In [304]: stacked.to_panel()
```

```
Out[304]:
```

```
<class 'pandas.core.panel.Panel'>
```

```
Dimensions: 6 (items) x 3 (major) x 4 (minor)
```

```
Items: Open to Adj Close
```

```
Major axis: 2012-05-30 00:00:00 to 2012-06-01 00:00:00
```

```
Minor axis: AAPL to MSFT
```

# Data Loading, Storage, and File Formats

The tools in this book are of little use if you can't easily import and export data in Python. I'm going to be focused on input and output with pandas objects, though there are of course numerous tools in other libraries to aid in this process. NumPy, for example, features low-level but extremely fast binary data loading and storage, including support for memory-mapped array. See [Chapter 12](#) for more on those.

Input and output typically falls into a few main categories: reading text files and other more efficient on-disk formats, loading data from databases, and interacting with network sources like web APIs.

## Reading and Writing Data in Text Format

Python has become a beloved language for text and file munging due to its simple syntax for interacting with files, intuitive data structures, and convenient features like tuple packing and unpacking.

pandas features a number of functions for reading tabular data as a DataFrame object. [Table 6-1](#) has a summary of all of them, though `read_csv` and `read_table` are likely the ones you'll use the most.

*Table 6-1. Parsing functions in pandas*

Function	Description
<code>read_csv</code>	Load delimited data from a file, URL, or file-like object. Use comma as default delimiter
<code>read_table</code>	Load delimited data from a file, URL, or file-like object. Use tab ('\\t') as default delimiter
<code>read_fwf</code>	Read data in fixed-width column format (that is, no delimiters)
<code>read_clipboard</code>	Version of <code>read_table</code> that reads data from the clipboard. Useful for converting tables from web pages

I'll give an overview of the mechanics of these functions, which are meant to convert text data into a DataFrame. The options for these functions fall into a few categories:

- Indexing: can treat one or more columns as the returned DataFrame, and whether to get column names from the file, the user, or not at all.
- Type inference and data conversion: this includes the user-defined value conversions and custom list of missing value markers.
- Datetime parsing: includes combining capability, including combining date and time information spread over multiple columns into a single column in the result.
- Iterating: support for iterating over chunks of very large files.
- Unclean data issues: skipping rows or a footer, comments, or other minor things like numeric data with thousands separated by commas.

*Type inference* is one of the more important features of these functions; that means you don't have to specify which columns are numeric, integer, boolean, or string. Handling dates and other custom types requires a bit more effort, though. Let's start with a small comma-separated (CSV) text file:

```
In [846]: !cat ch06/ex1.csv
a,b,c,d,message
1,2,3,4,hello
5,6,7,8,world
9,10,11,12,foo
```

Since this is comma-delimited, we can use `read_csv` to read it into a DataFrame:

```
In [847]: df = pd.read_csv('ch06/ex1.csv')

In [848]: df
Out[848]:
   a    b    c    d message
0  1    2    3    4    hello
1  5    6    7    8    world
2  9   10   11   12      foo
```

We could also have used `read_table` and specifying the delimiter:

```
In [849]: pd.read_table('ch06/ex1.csv', sep=',')
Out[849]:
   a    b    c    d message
0  1    2    3    4    hello
1  5    6    7    8    world
2  9   10   11   12      foo
```



Here I used the Unix `cat` shell command to print the raw contents of the file to the screen. If you're on Windows, you can use `type` instead of `cat` to achieve the same effect.

A file will not always have a header row. Consider this file:

```
In [850]: !cat ch06/ex2.csv
1,2,3,4,hello
5,6,7,8,world
9,10,11,12,foo
```

To read this in, you have a couple of options. You can allow pandas to assign default column names, or you can specify names yourself:

```
In [851]: pd.read_csv('ch06/ex2.csv', header=None)
Out[851]:
   X.1  X.2  X.3  X.4    X.5
0     1     2     3     4  hello
1     5     6     7     8  world
2     9    10    11    12    foo

In [852]: pd.read_csv('ch06/ex2.csv', names=['a', 'b', 'c', 'd', 'message'])
Out[852]:
      a    b    c    d message
0  1.0  2.0  3.0  4.0    hello
1  5.0  6.0  7.0  8.0    world
2  9.0 10.0 11.0 12.0     foo
```

Suppose you wanted the `message` column to be the index of the returned DataFrame. You can either indicate you want the column at index 4 or named '`message`' using the `index_col` argument:

```
In [853]: names = ['a', 'b', 'c', 'd', 'message']

In [854]: pd.read_csv('ch06/ex2.csv', names=names, index_col='message')
Out[854]:
      a    b    c    d
message
hello  1.0  2.0  3.0  4.0
world  5.0  6.0  7.0  8.0
foo    9.0 10.0 11.0 12.0
```

In the event that you want to form a hierarchical index from multiple columns, just pass a list of column numbers or names:

```
In [855]: !cat ch06/csv_mindex.csv
key1,key2,value1,value2
one,a,1,2
one,b,3,4
one,c,5,6
one,d,7,8
two,a,9,10
two,b,11,12
two,c,13,14
two,d,15,16
```

```
In [856]: parsed = pd.read_csv('ch06/csv_mindex.csv', index_col=['key1', 'key2'])

In [857]: parsed
Out[857]:
```

```

      value1  value2
key1 key2
one   a      1      2
      b      3      4
      c      5      6
      d      7      8
two   a      9     10
      b     11     12
      c     13     14
      d     15     16

```

In some cases, a table might not have a fixed delimiter, using whitespace or some other pattern to separate fields. In these cases, you can pass a regular expression as a delimiter for `read_table`. Consider a text file that looks like this:

```

In [858]: list(open('ch06/ex3.txt'))
Out[858]:
['          A          B          C\n',
 'aaa -0.264438 -1.026059 -0.619500\n',
 'bbb  0.927272  0.302904 -0.032399\n',
 'ccc -0.264273 -0.386314 -0.217601\n',
 'ddd -0.871858 -0.348382  1.100491\n']

```

While you could do some munging by hand, in this case fields are separated by a variable amount of whitespace. This can be expressed by the regular expression `\s+`, so we have then:

```

In [859]: result = pd.read_table('ch06/ex3.txt', sep='\s+')
In [860]: result
Out[860]:
          A          B          C
aaa -0.264438 -1.026059 -0.619500
bbb  0.927272  0.302904 -0.032399
ccc -0.264273 -0.386314 -0.217601
ddd -0.871858 -0.348382  1.100491

```

Because there was one fewer column name than the number of data rows, `read_table` infers that the first column should be the DataFrame's index in this special case.

The parser functions have many additional arguments to help you handle the wide variety of exception file formats that occur (see [Table 6-2](#)). For example, you can skip the first, third, and fourth rows of a file with `skiprows`:

```

In [861]: !cat ch06/ex4.csv
# hey!
a,b,c,d,message
# just wanted to make things more difficult for you
# who reads CSV files with computers, anyway?
1,2,3,4,hello
5,6,7,8,world
9,10,11,12,foo
In [862]: pd.read_csv('ch06/ex4.csv', skiprows=[0, 2, 3])
Out[862]:
     a     b     c     d message

```

```
0 1 2 3 4 hello
1 5 6 7 8 world
2 9 10 11 12 foo
```

Handling missing values is an important and frequently nuanced part of the file parsing process. Missing data is usually either not present (empty string) or marked by some *sentinel* value. By default, pandas uses a set of commonly occurring sentinels, such as NA, -1.#IND, and NULL:

```
In [863]: !cat ch06/ex5.csv
something,a,b,c,d,message
one,1,2,3,4,NA
two,5,6,,8,world
three,9,10,11,12,foo
In [864]: result = pd.read_csv('ch06/ex5.csv')

In [865]: result
Out[865]:
   something    a    b    c    d message
0      one    1    2    3    4      NaN
1     two    5    6  NaN    8    world
2   three    9   10   11   12      foo

In [866]: pd.isnull(result)
Out[866]:
   something      a      b      c      d message
0      False  False  False  False  False    True
1      False  False  False   True  False  False
2      False  False  False  False  False  False
```

The na\_values option can take either a list or set of strings to consider missing values:

```
In [867]: result = pd.read_csv('ch06/ex5.csv', na_values=['NULL'])

In [868]: result
Out[868]:
   something    a    b    c    d message
0      one    1    2    3    4      NaN
1     two    5    6  NaN    8    world
2   three    9   10   11   12      foo
```

Different NA sentinels can be specified for each column in a dict:

```
In [869]: sentinels = {'message': ['foo', 'NA'], 'something': ['two']}
In [870]: pd.read_csv('ch06/ex5.csv', na_values=sentinels)
Out[870]:
   something    a    b    c    d message
0      one    1    2    3    4      NaN
1      NaN    5    6  NaN    8    world
2   three    9   10   11   12      NaN
```

Table 6-2. `read_csv` /`read_table` function arguments

Argument	Description
<code>path</code>	String indicating filesystem location, URL, or file-like object
<code>sep</code> or <code>delimiter</code>	Character sequence or regular expression to use to split fields in each row
<code>header</code>	Row number to use as column names. Defaults to 0 (first row), but should be <code>None</code> if there is no header row
<code>index_col</code>	Column numbers or names to use as the row index in the result. Can be a single name/number or a list of them for a hierarchical index
<code>names</code>	List of column names for result, combine with <code>header=None</code>
<code>skiprows</code>	Number of rows at beginning of file to ignore or list of row numbers (starting from 0) to skip
<code>na_values</code>	Sequence of values to replace with <code>NA</code>
<code>comment</code>	Character or characters to split comments off the end of lines
<code>parse_dates</code>	Attempt to parse data to datetime; <code>False</code> by default. If <code>True</code> , will attempt to parse all columns. Otherwise can specify a list of column numbers or name to parse. If element of list is tuple or list, will combine multiple columns together and parse to date (for example if date/time split across two columns)
<code>keep_date_col</code>	If joining columns to parse date, drop the joined columns. Default <code>True</code>
<code>converters</code>	Dict containing column number or name mapping to functions. For example <code>{'foo': f}</code> would apply the function <code>f</code> to all values in the 'foo' column
<code>dayfirst</code>	When parsing potentially ambiguous dates, treat as international format (e.g. <code>7/6/2012</code> -> June 7, 2012). Default <code>False</code>
<code>date_parser</code>	Function to use to parse dates
<code>nrows</code>	Number of rows to read from beginning of file
<code>iterator</code>	Return a <code>TextParser</code> object for reading file piecemeal
<code>chunksize</code>	For iteration, size of file chunks
<code>skip_footer</code>	Number of lines to ignore at end of file
<code>verbose</code>	Print various parser output information, like the number of missing values placed in non-numeric columns
<code>encoding</code>	Text encoding for unicode. For example ' <code>utf-8</code> ' for UTF-8 encoded text
<code>squeeze</code>	If the parsed data only contains one column return a <code>Series</code>
<code>thousands</code>	Separator for thousands, e.g. <code>,</code> or <code>.</code>

## Reading Text Files in Pieces

When processing very large files or figuring out the right set of arguments to correctly process a large file, you may only want to read in a small piece of a file or iterate through smaller chunks of the file.

```
In [871]: result = pd.read_csv('ch06/ex6.csv')
```

```
In [872]: result
Out[872]:
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 10000 entries, 0 to 9999
Data columns:
one      10000 non-null values
two      10000 non-null values
three    10000 non-null values
four     10000 non-null values
key      10000 non-null values
dtypes: float64(4), object(1)
```

If you want to only read out a small number of rows (avoiding reading the entire file), specify that with `nrows`:

```
In [873]: pd.read_csv('ch06/ex6.csv', nrows=5)
Out[873]:
   one      two      three      four key
0  0.467976 -0.038649 -0.295344 -1.824726 L
1 -0.358893  1.404453  0.704965 -0.200638 B
2 -0.501840  0.659254 -0.421691 -0.057688 G
3  0.204886  1.074134  1.388361 -0.982404 R
4  0.354628 -0.133116  0.283763 -0.837063 Q
```

To read out a file in pieces, specify a `chunksize` as a number of rows:

```
In [874]: chunker = pd.read_csv('ch06/ex6.csv', chunksize=1000)

In [875]: chunker
Out[875]: <pandas.io.parsers.TextParser at 0x8398150>
```

The `TextParser` object returned by `read_csv` allows you to iterate over the parts of the file according to the `chunksize`. For example, we can iterate over `ex6.csv`, aggregating the value counts in the 'key' column like so:

```
chunker = pd.read_csv('ch06/ex6.csv', chunksize=1000)

tot = Series([])
for piece in chunker:
    tot = tot.add(piece['key'].value_counts(), fill_value=0)

tot = tot.order(ascending=False)
```

We have then:

```
In [877]: tot[:10]
Out[877]:
E    368
X    364
L    346
O    343
Q    340
M    338
J    337
F    335
K    334
H    330
```

`TextParser` is also equipped with a `get_chunk` method which enables you to read pieces of an arbitrary size.

## Writing Data Out to Text Format

Data can also be exported to delimited format. Let's consider one of the CSV files read above:

```
In [878]: data = pd.read_csv('ch06/ex5.csv')

In [879]: data
Out[879]:
   something  a    b    c    d  message
0      one  1    2    3    4     NaN
1      two  5    6  NaN  8    world
2    three  9   10   11   12     foo
```

Using DataFrame's `to_csv` method, we can write the data out to a comma-separated file:

```
In [880]: data.to_csv('ch06/out.csv')

In [881]: !cat ch06/out.csv
,something,a,b,c,d,message
0,one,1,2,3.0,4,
1,two,5,6,,8,world
2,three,9,10,11.0,12,foo
```

Other delimiters can be used, of course (writing to `sys.stdout` so it just prints the text result):

```
In [882]: data.to_csv(sys.stdout, sep='|')
|something|a|b|c|d|message
0|one|1|2|3.0|4|
1|two|5|6||8|world
2|three|9|10|11.0|12|foo
```

Missing values appear as empty strings in the output. You might want to denote them by some other sentinel value:

```
In [883]: data.to_csv(sys.stdout, na_rep='NULL')
,something,a,b,c,d,message
0,one,1,2,3.0,4,NULL
1,two,5,6,NULL,8,world
2,three,9,10,11.0,12,foo
```

With no other options specified, both the row and column labels are written. Both of these can be disabled:

```
In [884]: data.to_csv(sys.stdout, index=False, header=False)
one,1,2,3.0,4,
two,5,6,,8,world
three,9,10,11.0,12,foo
```

You can also write only a subset of the columns, and in an order of your choosing:

```
In [885]: data.to_csv(sys.stdout, index=False, cols=['a', 'b', 'c'])
a,b,c
1,2,3.0
5,6,
9,10,11.0
```

Series also has a `to_csv` method:

```
In [886]: dates = pd.date_range('1/1/2000', periods=7)
In [887]: ts = Series(np.arange(7), index=dates)
In [888]: ts.to_csv('ch06/tseries.csv')
```

```
In [889]: !cat ch06/tseries.csv
2000-01-01 00:00:00,0
2000-01-02 00:00:00,1
2000-01-03 00:00:00,2
2000-01-04 00:00:00,3
2000-01-05 00:00:00,4
2000-01-06 00:00:00,5
2000-01-07 00:00:00,6
```

With a bit of wrangling (no header, first column as index), you can read a CSV version of a Series with `read_csv`, but there is also a `from_csv` convenience method that makes it a bit simpler:

```
In [890]: Series.from_csv('ch06/tseries.csv', parse_dates=True)
Out[890]:
2000-01-01    0
2000-01-02    1
2000-01-03    2
2000-01-04    3
2000-01-05    4
2000-01-06    5
2000-01-07    6
```

See the docstrings for `to_csv` and `from_csv` in IPython for more information.

## Manually Working with Delimited Formats

Most forms of tabular data can be loaded from disk using functions like `pandas.read_table`. In some cases, however, some manual processing may be necessary. It's not uncommon to receive a file with one or more malformed lines that trip up `read_table`. To illustrate the basic tools, consider a small CSV file:

```
In [891]: !cat ch06/ex7.csv
"a","b","c"
"1","2","3"
"1","2","3","4"
```

For any file with a single-character delimiter, you can use Python's built-in `csv` module. To use it, pass any open file or file-like object to `csv.reader`:

```

import csv
f = open('ch06/ex7.csv')

reader = csv.reader(f)

```

Iterating through the reader like a file yields tuples of values in each line with any quote characters removed:

```

In [893]: for line in reader:
    ....:     print line
['a', 'b', 'c']
['1', '2', '3']
['1', '2', '3', '4']

```

From there, it's up to you to do the wrangling necessary to put the data in the form that you need it. For example:

```

In [894]: lines = list(csv.reader(open('ch06/ex7.csv')))

In [895]: header, values = lines[0], lines[1:]

In [896]: data_dict = {h: v for h, v in zip(header, *values)})

In [897]: data_dict
Out[897]: {'a': ('1', '1'), 'b': ('2', '2'), 'c': ('3', '3')}

```

CSV files come in many different flavors. Defining a new format with a different delimiter, string quoting convention, or line terminator is done by defining a simple subclass of `csv.Dialect`:

```

class my_dialect(csv.Dialect):
    lineterminator = '\n'
    delimiter = ';'
    quotechar = "'"

reader = csv.reader(f, dialect=my_dialect)

```

Individual CSV dialect parameters can also be given as keywords to `csv.reader` without having to define a subclass:

```
reader = csv.reader(f, delimiter='|')
```

The possible options (attributes of `csv.Dialect`) and what they do can be found in [Table 6-3](#).

*Table 6-3. CSV dialect options*

Argument	Description
delimiter	One-character string to separate fields. Defaults to ','.
lineterminator	Line terminator for writing, defaults to '\r\n'. Reader ignores this and recognizes cross-platform line terminators.
quotechar	Quote character for fields with special characters (like a delimiter). Default is "'".
quoting	Quoting convention. Options include <code>csv.QUOTE_ALL</code> (quote all fields), <code>csv.QUOTE_MINIMAL</code> (only fields with special characters like the delimiter),

Argument	Description
	<code>csv.QUOTE_NONNUMERIC</code> , and <code>csv.QUOTE_NON</code> (no quoting). See Python's documentation for full details. Defaults to <code>QUOTE_MINIMAL</code> .
<code>skipinitialspace</code>	Ignore whitespace after each delimiter. Default <code>False</code> .
<code>doublequote</code>	How to handle quoting character inside a field. If <code>True</code> , it is doubled. See online documentation for full detail and behavior.
<code>escapechar</code>	String to escape the delimiter if <code>quoting</code> is set to <code>csv.QUOTE_NONE</code> . Disabled by default



For files with more complicated or fixed multicharacter delimiters, you will not be able to use the `csv` module. In those cases, you'll have to do the line splitting and other cleanup using string's `split` method or the regular expression method `re.split`.

To *write* delimited files manually, you can use `csv.writer`. It accepts an open, writable file object and the same dialect and format options as `csv.reader`:

```
with open('mydata.csv', 'w') as f:
    writer = csv.writer(f, dialect='my_dialect')
    writer.writerow(['one', 'two', 'three'])
    writer.writerow(['1', '2', '3'])
    writer.writerow(['4', '5', '6'])
    writer.writerow(['7', '8', '9'])
```

## JSON Data

JSON (short for JavaScript Object Notation) has become one of the standard formats for sending data by HTTP request between web browsers and other applications. It is a much more flexible data format than a tabular text form like CSV. Here is an example:

```
obj = """
{
    "name": "Wes",
    "places_lived": ["United States", "Spain", "Germany"],
    "pet": null,
    "siblings": [{"name": "Scott", "age": 25, "pet": "Zuko"},
                 {"name": "Katie", "age": 33, "pet": "Cisco"}]
}
"""
```

JSON is very nearly valid Python code with the exception of its null value `null` and some other nuances (such as disallowing trailing commas at the end of lists). The basic types are objects (dicts), arrays (lists), strings, numbers, booleans, and nulls. All of the keys in an object must be strings. There are several Python libraries for reading and writing JSON data. I'll use `json` here as it is built into the Python standard library. To convert a JSON string to Python form, use `json.loads`:

In [899]: `import json`

```
In [900]: result = json.loads(obj)

In [901]: result
Out[901]:
{u'name': u'Wes',
 u'pet': None,
 u'places_lived': [u'United States', u'Spain', u'Germany'],
 u'siblings': [{u'age': 25, u'name': u'Scott', u'pet': u'Zuko'},
   {u'age': 33, u'name': u'Katie', u'pet': u'Cisco'}]}
```

`json.dumps` on the other hand converts a Python object back to JSON:

```
In [902]: asjson = json.dumps(result)
```

How you convert a JSON object or list of objects to a DataFrame or some other data structure for analysis will be up to you. Conveniently, you can pass a list of JSON objects to the DataFrame constructor and select a subset of the data fields:

```
In [903]: siblings = DataFrame(result['siblings'], columns=['name', 'age'])

In [904]: siblings
Out[904]:
   name  age
0  Scott   25
1  Katie   33
```

For an extended example of reading and manipulating JSON data (including nested records), see the USDA Food Database example in the next chapter.



An effort is underway to add fast native JSON export (`to_json`) and decoding (`from_json`) to pandas. This was not ready at the time of writing.

## XML and HTML: Web Scraping

Python has many libraries for reading and writing data in the ubiquitous HTML and XML formats. `lxml` (<http://lxml.de>) is one that has consistently strong performance in parsing very large files. `lxml` has multiple programmer interfaces; first I'll show using `lxml.html` for HTML, then parse some XML using `lxml.objectify`.

Many websites make data available in HTML tables for viewing in a browser, but not downloadable as an easily machine-readable format like JSON, HTML, or XML. I noticed that this was the case with Yahoo! Finance's stock options data. If you aren't familiar with this data; options are derivative contracts giving you the right to buy (*call* option) or sell (*put* option) a company's stock at some particular price (the *strike*) between now and some fixed point in the future (the *expiry*). People trade both *call* and *put* options across many strikes and expiries; this data can all be found together in tables on Yahoo! Finance.

To get started, find the URL you want to extract data from, open it with `urllib2` and parse the stream with `lxml` like so:

```
from lxml.html import parse
from urllib2 import urlopen

parsed = parse(urlopen('http://finance.yahoo.com/q/op?s=AAPL+Options'))

doc = parsed.getroot()
```

Using this object, you can extract all HTML tags of a particular type, such as `table` tags containing the data of interest. As a simple motivating example, suppose you wanted to get a list of every URL linked to in the document; links are `a` tags in HTML. Using the document root's `findall` method along with an XPath (a means of expressing "queries" on the document):

```
In [906]: links = doc.findall('.//a')

In [907]: links[15:20]
Out[907]:
[<Element a at 0x6c488f0>,
 <Element a at 0x6c48950>,
 <Element a at 0x6c489b0>,
 <Element a at 0x6c48a10>,
 <Element a at 0x6c48a70>]
```

But these are objects representing HTML elements; to get the URL and link text you have to use each element's `get` method (for the URL) and `text_content` method (for the display text):

```
In [908]: lnk = links[28]

In [909]: lnk
Out[909]: <Element a at 0x6c48dd0>

In [910]: lnk.get('href')
Out[910]: 'http://biz.yahoo.com/special.html'

In [911]: lnk.text_content()
Out[911]: 'Special Editions'
```

Thus, getting a list of all URLs in the document is a matter of writing this list comprehension:

```
In [912]: urls = [lnk.get('href') for lnk in doc.findall('.//a')]

In [913]: urls[-10:]
Out[913]:
['http://info.yahoo.com/privacy/us/yahoo/finance/details.html',
 'http://info.yahoo.com/relevantads/',
 'http://docs.yahoo.com/info/terms/',
 'http://docs.yahoo.com/info/copyright/copyright.html',
 'http://help.yahoo.com/l/us/yahoo/finance/forms_index.html',
 'http://help.yahoo.com/l/us/yahoo/finance/quotes/fitadelay.html',
 'http://help.yahoo.com/l/us/yahoo/finance/quotes/fitadelay.html',
```

```
'http://www.capitaliq.com',
'http://www.csidata.com',
'http://www.morningstar.com/']
```

Now, finding the right tables in the document can be a matter of trial and error; some websites make it easier by giving a table of interest an `id` attribute. I determined that these were the two tables containing the call data and put data, respectively:

```
tables = doc.findall('.//table')
calls = tables[9]
puts = tables[13]
```

Each table has a header row followed by each of the data rows:

```
In [915]: rows = calls.findall('.//tr')
```

For the header as well as the data rows, we want to extract the text from each cell; in the case of the header these are `th` cells and `td` cells for the data:

```
def _unpack(row, kind='td'):
    elts = row.findall('.//%s' % kind)
    return [val.text_content() for val in elts]
```

Thus, we obtain:

```
In [917]: _unpack(rows[0], kind='th')
Out[917]: ['Strike', 'Symbol', 'Last', 'Chg', 'Bid', 'Ask', 'Vol', 'Open Int']
```

```
In [918]: _unpack(rows[1], kind='td')
Out[918]:
['295.00',
'AAPL120818C00295000',
'310.40',
'0.00',
'289.80',
'290.80',
'1',
'169']
```

Now, it's a matter of combining all of these steps together to convert this data into a DataFrame. Since the numerical data is still in string format, we want to convert some, but perhaps not all of the columns to floating point format. You could do this by hand, but, luckily, pandas has a class `TextParser` that is used internally in the `read_csv` and other parsing functions to do the appropriate automatic type conversion:

```
from pandas.io.parsers import TextParser

def parse_options_data(table):
    rows = table.findall('.//tr')
    header = _unpack(rows[0], kind='th')
    data = [_unpack(r) for r in rows[1:]]
    return TextParser(data, names=header).get_chunk()
```

Finally, we invoke this parsing function on the lxml table objects and get DataFrame results:

```
In [920]: call_data = parse_options_data(calls)
```

```
In [921]: put_data = parse_options_data(puts)
```

```
In [922]: call_data[:10]
```

```
Out[922]:
```

	Strike	Symbol	Last	Chg	Bid	Ask	Vol	Open	Int
0	295	AAPL120818C00295000	310.40	0.0	289.80	290.80	1	169	
1	300	AAPL120818C00300000	277.10	1.7	284.80	285.60	2	478	
2	305	AAPL120818C00305000	300.97	0.0	279.80	280.80	10	316	
3	310	AAPL120818C00310000	267.05	0.0	274.80	275.65	6	239	
4	315	AAPL120818C00315000	296.54	0.0	269.80	270.80	22	88	
5	320	AAPL120818C00320000	291.63	0.0	264.80	265.80	96	173	
6	325	AAPL120818C00325000	261.34	0.0	259.80	260.80	N/A	108	
7	330	AAPL120818C00330000	230.25	0.0	254.80	255.80	N/A	21	
8	335	AAPL120818C00335000	266.03	0.0	249.80	250.65	4	46	
9	340	AAPL120818C00340000	272.58	0.0	244.80	245.80	4	30	

## Parsing XML with lxml.objectify

XML (extensible markup language) is another common structured data format supporting hierarchical, nested data with metadata. The files that generate the book you are reading actually form a series of large XML documents.

Above, I showed the `lxml` library and its `lxml.html` interface. Here I show an alternate interface that's convenient for XML data, `lxml.objectify`.

The New York Metropolitan Transportation Authority (MTA) publishes a number of data series about its bus and train services (<http://www.mta.info/developers/download.html>). Here we'll look at the performance data which is contained in a set of XML files. Each train or bus service has a different file (like `Performance_MNR.xml` for the Metro-North Railroad) containing monthly data as a series of XML records that look like this:

```
<INDICATOR>
  <INDICATOR_SEQ>373889</INDICATOR_SEQ>
  <PARENT_SEQ></PARENT_SEQ>
  <AGENCY_NAME>Metro-North Railroad</AGENCY_NAME>
  <INDICATOR_NAME>Escalator Availability</INDICATOR_NAME>
  <DESCRIPTION>Percent of the time that escalators are operational
  systemwide. The availability rate is based on physical observations performed
  the morning of regular business days only. This is a new indicator the agency
  began reporting in 2009.</DESCRIPTION>
  <PERIOD_YEAR>2011</PERIOD_YEAR>
  <PERIOD_MONTH>12</PERIOD_MONTH>
  <CATEGORY>Service Indicators</CATEGORY>
  <FREQUENCY>M</FREQUENCY>
  <DESIRED_CHANGE>U</DESIRED_CHANGE>
  <INDICATOR_UNIT>%</INDICATOR_UNIT>
  <DECIMAL_PLACES>1</DECIMAL_PLACES>
  <YTD_TARGET>97.00</YTD_TARGET>
  <YTD_ACTUAL></YTD_ACTUAL>
  <MONTHLY_TARGET>97.00</MONTHLY_TARGET>
  <MONTHLY_ACTUAL></MONTHLY_ACTUAL>
</INDICATOR>
```

Using `lxml.objectify`, we parse the file and get a reference to the root node of the XML file with `getroot`:

```
from lxml import objectify

path = 'Performance_MNR.xml'
parsed = objectify.parse(open(path))
root = parsed.getroot()
```

`root.INDICATOR` return a generator yielding each `<INDICATOR>` XML element. For each record, we can populate a dict of tag names (like `YTD_ACTUAL`) to data values (excluding a few tags):

```
data = []

skip_fields = ['PARENT_SEQ', 'INDICATOR_SEQ',
               'DESIRED_CHANGE', 'DECIMAL_PLACES']

for elt in root.INDICATOR:
    el_data = {}
    for child in elt.getchildren():
        if child.tag in skip_fields:
            continue
        el_data[child.tag] = child.pyval
    data.append(el_data)
```

Lastly, convert this list of dicts into a DataFrame:

```
In [927]: perf = DataFrame(data)
```

```
In [928]: perf
Out[928]:
Empty DataFrame
Columns: array([], dtype=int64)
Index: array([], dtype=int64)
```

XML data can get much more complicated than this example. Each tag can have metadata, too. Consider an HTML link tag which is also valid XML:

```
from StringIO import StringIO
tag = '<a href="http://www.google.com">Google</a>'

root = objectify.parse(StringIO(tag)).getroot()
```

You can now access any of the fields (like `href`) in the tag or the link text:

```
In [930]: root
Out[930]: <Element a at 0x88bd4b0>

In [931]: root.get('href')
Out[931]: 'http://www.google.com'

In [932]: root.text
Out[932]: 'Google'
```

## Binary Data Formats

One of the easiest ways to store data efficiently in binary format is using Python's built-in `pickle` serialization. Conveniently, pandas objects all have a `save` method which writes the data to disk as a pickle:

```
In [933]: frame = pd.read_csv('ch06/ex1.csv')
```

```
In [934]: frame
```

```
Out[934]:
```

	a	b	c	d	message
0	1	2	3	4	hello
1	5	6	7	8	world
2	9	10	11	12	foo

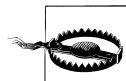
```
In [935]: frame.save('ch06/frame_pickle')
```

You read the data back into Python with `pandas.load`, another pickle convenience function:

```
In [936]: pd.load('ch06/frame_pickle')
```

```
Out[936]:
```

	a	b	c	d	message
0	1	2	3	4	hello
1	5	6	7	8	world
2	9	10	11	12	foo



`pickle` is only recommended as a short-term storage format. The problem is that it is hard to guarantee that the format will be stable over time; an object pickled today may not unpickle with a later version of a library. I have made every effort to ensure that this does not occur with pandas, but at some point in the future it may be necessary to “break” the pickle format.

## Using HDF5 Format

There are a number of tools that facilitate efficiently reading and writing large amounts of scientific data in binary format on disk. A popular industry-grade library for this is HDF5, which is a C library with interfaces in many other languages like Java, Python, and MATLAB. The “HDF” in HDF5 stands for *hierarchical data format*. Each HDF5 file contains an internal file system-like node structure enabling you to store multiple datasets and supporting metadata. Compared with simpler formats, HDF5 supports on-the-fly compression with a variety of compressors, enabling data with repeated patterns to be stored more efficiently. For very large datasets that don’t fit into memory, HDF5 is a good choice as you can efficiently read and write small sections of much larger arrays.

There are not one but two interfaces to the HDF5 library in Python, PyTables and h5py, each of which takes a different approach to the problem. h5py provides a direct, but high-level interface to the HDF5 API, while PyTables abstracts many of the details of

HDF5 to provide multiple flexible data containers, table indexing, querying capability, and some support for out-of-core computations.

pandas has a minimal dict-like `HDFStore` class, which uses PyTables to store pandas objects:

```
In [937]: store = pd.HDFStore('mydata.h5')

In [938]: store['obj1'] = frame

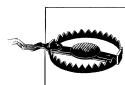
In [939]: store['obj1_col'] = frame['a']

In [940]: store
Out[940]:
<class 'pandas.io.pytables.HDFStore'>
File path: mydata.h5
obj1           DataFrame
obj1_col        Series
```

Objects contained in the HDF5 file can be retrieved in a dict-like fashion:

```
In [941]: store['obj1']
Out[941]:
   a   b   c   d  message
0  1   2   3   4    hello
1  5   6   7   8   world
2  9  10  11  12     foo
```

If you work with huge quantities of data, I would encourage you to explore PyTables and h5py to see how they can suit your needs. Since many data analysis problems are IO-bound (rather than CPU-bound), using a tool like HDF5 can massively accelerate your applications.



HDF5 is *not* a database. It is best suited for write-once, read-many datasets. While data can be added to a file at any time, if multiple writers do so simultaneously, the file can become corrupted.

## Reading Microsoft Excel Files

pandas also supports reading tabular data stored in Excel 2003 (and higher) files using the `ExcelFile` class. Internally `ExcelFile` uses the `xlrd` and `openpyxl` packages, so you may have to install them first. To use `ExcelFile`, create an instance by passing a path to an `xls` or `xlsx` file:

```
xls_file = pd.ExcelFile('data.xls')
```

Data stored in a sheet can then be read into `DataFrame` using `parse`:

```
table = xls_file.parse('Sheet1')
```

## Interacting with HTML and Web APIs

Many websites have public APIs providing data feeds via JSON or some other format. There are a number of ways to access these APIs from Python; one easy-to-use method that I recommend is the `requests` package (<http://docs.python-requests.org>). To search for the words “python pandas” on Twitter, we can make an HTTP GET request like so:

```
In [944]: import requests  
In [945]: url = 'http://search.twitter.com/search.json?q=python%20pandas'  
In [946]: resp = requests.get(url)  
In [947]: resp  
Out[947]: <Response [200]>
```

The Response object’s `text` attribute contains the content of the GET query. Many web APIs will return a JSON string that must be loaded into a Python object:

```
In [948]: import json  
In [949]: data = json.loads(resp.text)  
In [950]: data.keys()  
Out[950]:  
[u'next_page',  
 u'completed_in',  
 u'max_id_str',  
 u'since_id_str',  
 u'refresh_url',  
 u'results',  
 u'since_id',  
 u'results_per_page',  
 u'query',  
 u'max_id',  
 u'page']
```

The `results` field in the response contains a list of tweets, each of which is represented as a Python dict that looks like:

```
{u'created_at': u'Mon, 25 Jun 2012 17:50:33 +0000',  
 u'from_user': u'wesmckinn',  
 u'from_user_id': 115494880,  
 u'from_user_id_str': u'115494880',  
 u'from_user_name': u'Wes McKinney',  
 u'geo': None,  
 u'id': 217313849177686018,  
 u'id_str': u'217313849177686018',  
 u'iso_language_code': u'pt',  
 u'metadata': {u'result_type': u'recent'},  
 u'source': u'<a href="http://twitter.com/">web</a>',  
 u'text': u'Lunchtime pandas-fu http://t.co/SI70xZZQ #pydata',  
 u'to_user': None,  
 u'to_user_id': 0,
```

```
u'to_user_id_str': u'0',
u'to_user_name': None}
```

We can then make a list of the tweet fields of interest then pass the results list to DataFrame:

```
In [951]: tweet_fields = ['created_at', 'from_user', 'id', 'text']

In [952]: tweets = DataFrame(data['results'], columns=tweet_fields)

In [953]: tweets
Out[953]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 15 entries, 0 to 14
Data columns:
created_at    15 non-null values
from_user     15 non-null values
id            15 non-null values
text          15 non-null values
dtypes: int64(1), object(3)
```

Each row in the DataFrame now has the extracted data from each tweet:

```
In [121]: tweets.ix[7]
Out[121]:
created at           Thu, 23 Jul 2012 09:54:00 +0000
from_user            deblike
id                  227419585803059201
text      pandas: powerful Python data analysis toolkit
Name: 7
```

With a bit of elbow grease, you can create some higher-level interfaces to common web APIs that return DataFrame objects for easy analysis.

## Interacting with Databases

In many applications data rarely comes from text files, that being a fairly inefficient way to store large amounts of data. SQL-based relational databases (such as SQL Server, PostgreSQL, and MySQL) are in wide use, and many alternative non-SQL (so-called *NoSQL*) databases have become quite popular. The choice of database is usually dependent on the performance, data integrity, and scalability needs of an application.

Loading data from SQL into a DataFrame is fairly straightforward, and pandas has some functions to simplify the process. As an example, I'll use an in-memory SQLite database using Python's built-in `sqlite3` driver:

```
import sqlite3

query = """
CREATE TABLE test
(a VARCHAR(20), b VARCHAR(20),
 c REAL,          d INTEGER
);"""
```

```
con = sqlite3.connect(':memory:')
con.execute(query)
con.commit()
```

Then, insert a few rows of data:

```
data = [('Atlanta', 'Georgia', 1.25, 6),
        ('Tallahassee', 'Florida', 2.6, 3),
        ('Sacramento', 'California', 1.7, 5)]
stmt = "INSERT INTO test VALUES(?, ?, ?, ?)"

con.executemany(stmt, data)
con.commit()
```

Most Python SQL drivers (PyODBC, psycopg2, MySQLdb, pymssql, etc.) return a list of tuples when selecting data from a table:

```
In [956]: cursor = con.execute('select * from test')
```

```
In [957]: rows = cursor.fetchall()
```

```
In [958]: rows
```

```
Out[958]:
```

```
[('Atlanta', u'Georgia', 1.25, 6),
 ('Tallahassee', u'Florida', 2.6, 3),
 ('Sacramento', u'California', 1.7, 5)]
```

You can pass the list of tuples to the DataFrame constructor, but you also need the column names, contained in the cursor's `description` attribute:

```
In [959]: cursor.description
```

```
Out[959]:
```

```
(('a', None, None, None, None, None, None),
 ('b', None, None, None, None, None, None),
 ('c', None, None, None, None, None, None),
 ('d', None, None, None, None, None, None))
```

```
In [960]: DataFrame(rows, columns=zip(*cursor.description)[0])
```

```
Out[960]:
```

	a	b	c	d
0	Atlanta	Georgia	1.25	6
1	Tallahassee	Florida	2.60	3
2	Sacramento	California	1.70	5

This is quite a bit of munging that you'd rather not repeat each time you query the database. pandas has a `read_frame` function in its `pandas.io.sql` module that simplifies the process. Just pass the select statement and the connection object:

```
In [961]: import pandas.io.sql as sql
```

```
In [962]: sql.read_frame('select * from test', con)
```

```
Out[962]:
```

	a	b	c	d
0	Atlanta	Georgia	1.25	6
1	Tallahassee	Florida	2.60	3
2	Sacramento	California	1.70	5

## Storing and Loading Data in MongoDB

NoSQL databases take many different forms. Some are simple dict-like key-value stores like BerkeleyDB or Tokyo Cabinet, while others are document-based, with a dict-like object being the basic unit of storage. I've chosen MongoDB (<http://mongodb.org>) for my example. I started a MongoDB instance locally on my machine, and connect to it on the default port using `pymongo`, the official driver for MongoDB:

```
import pymongo
con = pymongo.Connection('localhost', port=27017)
```

Documents stored in MongoDB are found in collections inside databases. Each running instance of the MongoDB server can have multiple databases, and each database can have multiple collections. Suppose I wanted to store the Twitter API data from earlier in the chapter. First, I can access the (currently empty) tweets collection:

```
tweets = con.db.tweets
```

Then, I load the list of tweets and write each of them to the collection using `tweets.save` (which writes the Python dict to MongoDB):

```
import requests, json
url = 'http://search.twitter.com/search.json?q=python%20pandas'
data = json.loads(requests.get(url).text)

for tweet in data['results']:
    tweets.save(tweet)
```

Now, if I wanted to get all of my tweets (if any) from the collection, I can query the collection with the following syntax:

```
cursor = tweets.find({'from_user': 'wesmckinn'})
```

The cursor returned is an iterator that yields each document as a dict. As above I can convert this into a DataFrame, optionally extracting a subset of the data fields in each tweet:

```
tweet_fields = ['created_at', 'from_user', 'id', 'text']
result = DataFrame(list(cursor), columns=tweet_fields)
```

# Data Wrangling: Clean, Transform, Merge, Reshape

Much of the programming work in data analysis and modeling is spent on data preparation: loading, cleaning, transforming, and rearranging. Sometimes the way that data is stored in files or databases is not the way you need it for a data processing application. Many people choose to do ad hoc processing of data from one form to another using a general purpose programming, like Python, Perl, R, or Java, or UNIX text processing tools like sed or awk. Fortunately, pandas along with the Python standard library provide you with a high-level, flexible, and high-performance set of core manipulations and algorithms to enable you to wrangle data into the right form without much trouble.

If you identify a type of data manipulation that isn't anywhere in this book or elsewhere in the pandas library, feel free to suggest it on the mailing list or GitHub site. Indeed, much of the design and implementation of pandas has been driven by the needs of real world applications.

## Combining and Merging Data Sets

Data contained in pandas objects can be combined together in a number of built-in ways:

- `pandas.merge` connects rows in DataFrames based on one or more keys. This will be familiar to users of SQL or other relational databases, as it implements database *join* operations.
- `pandas.concat` glues or stacks together objects along an axis.
- `combine_first` instance method enables splicing together overlapping data to fill in missing values in one object with values from another.

I will address each of these and give a number of examples. They'll be utilized in examples throughout the rest of the book.

## Database-style DataFrame Merges

*Merge* or *join* operations combine data sets by linking rows using one or more *keys*. These operations are central to relational databases. The `merge` function in pandas is the main entry point for using these algorithms on your data.

Let's start with a simple example:

```
In [15]: df1 = DataFrame({'key': ['b', 'b', 'a', 'c', 'a', 'a', 'b'],
....:                   'data1': range(7)})

In [16]: df2 = DataFrame({'key': ['a', 'b', 'd'],
....:                   'data2': range(3)})

In [17]: df1
Out[17]:
   data1  key
0      0    b
1      1    b
2      2    a
3      3    c
4      4    a
5      5    a
6      6    b

In [18]: df2
Out[18]:
   data2  key
0      0    a
1      1    b
2      2    d
```

This is an example of a *many-to-one* merge situation; the data in `df1` has multiple rows labeled `a` and `b`, whereas `df2` has only one row for each value in the `key` column. Calling `merge` with these objects we obtain:

```
In [19]: pd.merge(df1, df2)
Out[19]:
   data1  key  data2
0      2    a      0
1      4    a      0
2      5    a      0
3      0    b      1
4      1    b      1
5      6    b      1
```

Note that I didn't specify which column to join on. If not specified, `merge` uses the overlapping column names as the keys. It's a good practice to specify explicitly, though:

```
In [20]: pd.merge(df1, df2, on='key')
Out[20]:
   data1  key  data2
0      2    a      0
1      4    a      0
2      5    a      0
3      0    b      1
4      1    b      1
5      6    b      1
```

If the column names are different in each object, you can specify them separately:

```
In [21]: df3 = DataFrame({'lkey': ['b', 'b', 'a', 'c', 'a', 'a', 'b'],
....:                   'data1': range(7)})
```

```
In [22]: df4 = DataFrame({'rkey': ['a', 'b', 'd'],
....:                      'data2': range(3)})

In [23]: pd.merge(df3, df4, left_on='lkey', right_on='rkey')
Out[23]:
   data1  lkey  data2  rkey
0      2     a      0     a
1      4     a      0     a
2      5     a      0     a
3      0     b      1     b
4      1     b      1     b
5      6     b      1     b
```

You probably noticed that the 'c' and 'd' values and associated data are missing from the result. By default `merge` does an '`inner`' join; the keys in the result are the intersection. Other possible options are '`left`', '`right`', and '`outer`'. The outer join takes the union of the keys, combining the effect of applying both left and right joins:

```
In [24]: pd.merge(df1, df2, how='outer')
Out[24]:
   data1  key  data2
0      2     a      0
1      4     a      0
2      5     a      0
3      0     b      1
4      1     b      1
5      6     b      1
6      3     c    NaN
7     NaN     d      2
```

*Many-to-many* merges have well-defined though not necessarily intuitive behavior. Here's an example:

```
In [25]: df1 = DataFrame({'key': ['b', 'b', 'a', 'c', 'a', 'b'],
....:                      'data1': range(6)})

In [26]: df2 = DataFrame({'key': ['a', 'b', 'a', 'b', 'd'],
....:                      'data2': range(5)})

In [27]: df1
Out[27]:
   data1  key
0      0     b
1      1     b
2      2     a
3      3     c
4      4     a
5      5     b

In [28]: df2
Out[28]:
   data2  key
0      0     a
1      1     b
2      2     a
3      3     b
4      4     d

In [29]: pd.merge(df1, df2, on='key', how='left')
Out[29]:
   data1  key  data2
0      2     a      0
1      2     a      2
```

```

2      4   a    0
3      4   a    2
4      0   b    1
5      0   b    3
6      1   b    1
7      1   b    3
8      5   b    1
9      5   b    3
10     3   c    NaN

```

Many-to-many joins form the Cartesian product of the rows. Since there were 3 'b' rows in the left DataFrame and 2 in the right one, there are 6 'b' rows in the result. The join method only affects the distinct key values appearing in the result:

```

In [30]: pd.merge(df1, df2, how='inner')
Out[30]:
   data1  key  data2
0      2   a    0
1      2   a    2
2      4   a    0
3      4   a    2
4      0   b    1
5      0   b    3
6      1   b    1
7      1   b    3
8      5   b    1
9      5   b    3

```

To merge with multiple keys, pass a list of column names:

```

In [31]: left = DataFrame({'key1': ['foo', 'foo', 'bar'],
....:                   'key2': ['one', 'two', 'one'],
....:                   'lval': [1, 2, 3]})

In [32]: right = DataFrame({'key1': ['foo', 'foo', 'bar', 'bar'],
....:                   'key2': ['one', 'one', 'one', 'two'],
....:                   'rval': [4, 5, 6, 7]})

In [33]: pd.merge(left, right, on=['key1', 'key2'], how='outer')
Out[33]:
   key1  key2  lval  rval
0  bar   one     3     6
1  bar   two    NaN     7
2  foo   one     1     4
3  foo   one     1     5
4  foo   two     2    NaN

```

To determine which key combinations will appear in the result depending on the choice of merge method, think of the multiple keys as forming an array of tuples to be used as a single join key (even though it's not actually implemented that way).



When joining columns-on-columns, the indexes on the passed DataFrame objects are discarded.

A last issue to consider in merge operations is the treatment of overlapping column names. While you can address the overlap manually (see the later section on renaming axis labels), `merge` has a `suffixes` option for specifying strings to append to overlapping names in the left and right DataFrame objects:

```
In [34]: pd.merge(left, right, on='key1')
Out[34]:
   key1  key2_x  lval  key2_y  rval
0  bar      one     3    one     6
1  bar      one     3   two     7
2  foo      one     1    one     4
3  foo      one     1    one     5
4  foo     two     2    one     4
5  foo     two     2    one     5

In [35]: pd.merge(left, right, on='key1', suffixes=('_left', '_right'))
Out[35]:
   key1  key2_left  lval  key2_right  rval
0  bar      one     3        one     6
1  bar      one     3       two     7
2  foo      one     1        one     4
3  foo      one     1        one     5
4  foo     two     2        one     4
5  foo     two     2        one     5
```

See [Table 7-1](#) for an argument reference on `merge`. Joining on index is the subject of the next section.

*Table 7-1. merge function arguments*

Argument	Description
left	DataFrame to be merged on the left side
right	DataFrame to be merged on the right side
how	One of 'inner', 'outer', 'left' or 'right'. 'inner' by default
on	Column names to join on. Must be found in both DataFrame objects. If not specified and no other join keys given, will use the intersection of the column names in <code>left</code> and <code>right</code> as the join keys
left_on	Columns in <code>left</code> DataFrame to use as join keys
right_on	Analogous to <code>left_on</code> for <code>left</code> DataFrame
left_index	Use row index in <code>left</code> as its join key (or keys, if a MultiIndex)
right_index	Analogous to <code>left_index</code>
sort	Sort merged data lexicographically by join keys; True by default. Disable to get better performance in some cases on large datasets
suffixes	Tuple of string values to append to column names in case of overlap; defaults to ('_x', '_y'). For example, if 'data' in both DataFrame objects, would appear as 'data_x' and 'data_y' in result
copy	If False, avoid copying data into resulting data structure in some exceptional cases. By default always copies

## Merging on Index

In some cases, the merge key or keys in a DataFrame will be found in its index. In this case, you can pass `left_index=True` or `right_index=True` (or both) to indicate that the index should be used as the merge key:

```
In [36]: left1 = DataFrame({'key': ['a', 'b', 'a', 'a', 'b', 'c'],
....:                      'value': range(6)})
```

```
In [37]: right1 = DataFrame({'group_val': [3.5, 7]}, index=['a', 'b'])
```

```
In [38]: left1           In [39]: right1
Out[38]:          Out[39]:
   key  value        group_val
0    a     0            a      3.5
1    b     1            b      7.0
2    a     2
3    a     3
4    b     4
5    c     5
```

```
In [40]: pd.merge(left1, right1, left_on='key', right_index=True)
```

```
Out[40]:
   key  value  group_val
0    a     0      3.5
1    a     2      3.5
2    a     3      3.5
3    b     1      7.0
4    b     4      7.0
```

Since the default merge method is to intersect the join keys, you can instead form the union of them with an outer join:

```
In [41]: pd.merge(left1, right1, left_on='key', right_index=True, how='outer')
Out[41]:
   key  value  group_val
0    a     0      3.5
1    a     2      3.5
2    a     3      3.5
3    b     1      7.0
4    b     4      7.0
5    c     5      NaN
```

With hierarchically-indexed data, things are a bit more complicated:

```
In [42]: lefth = DataFrame({'key1': ['Ohio', 'Ohio', 'Ohio', 'Nevada', 'Nevada'],
....:                      'key2': [2000, 2001, 2002, 2001, 2002],
....:                      'data': np.arange(5.)})
```

```
In [43]: righth = DataFrame(np.arange(12).reshape((6, 2)),
....:                      index=[['Nevada', 'Nevada', 'Ohio', 'Ohio', 'Ohio', 'Ohio'],
....:                             [2001, 2000, 2000, 2000, 2001, 2002]],
....:                      columns=['event1', 'event2'])
```

```
In [44]: lefth           In [45]: righth
Out[44]:          Out[45]:
```

```

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

```

In this case, you have to indicate multiple columns to merge on as a list (pay attention to the handling of duplicate index values):

```
In [46]: pd.merge(left, right, left_on=['key1', 'key2'], right_index=True)
```

```
Out[46]:
```

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

```
In [47]: pd.merge(left, right, left_on=['key1', 'key2'],
....:                 right_index=True, how='outer')
```

```
Out[47]:
```

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

Using the indexes of both sides of the merge is also not an issue:

```
In [48]: left2 = DataFrame([[1., 2.], [3., 4.], [5., 6.]], index=['a', 'c', 'e'],
....:                      columns=['Ohio', 'Nevada'])
```

```
In [49]: right2 = DataFrame([[7., 8.], [9., 10.], [11., 12.], [13, 14.]],
....:                      index=['b', 'c', 'd', 'e'], columns=['Missouri', 'Alabama'])
```

```
In [50]: left2
```

```
Out[50]:
```

	Ohio	Nevada	Missouri	Alabama	
a	1	2	b	7	8
c	3	4	c	9	10
e	5	6	d	11	12
			e	13	14

```
In [51]: right2
```

```
Out[51]:
```

```
In [52]: pd.merge(left2, right2, how='outer', left_index=True, right_index=True)
```

```
Out[52]:
```

	Ohio	Nevada	Missouri	Alabama
a	1	2	NaN	NaN
b	NaN	NaN	7	8
c	3	4	9	10
d	NaN	NaN	11	12
e	5	6	13	14

DataFrame has a more convenient `join` instance for merging by index. It can also be used to combine together many DataFrame objects having the same or similar indexes but non-overlapping columns. In the prior example, we could have written:

```
In [53]: left2.join(right2, how='outer')
Out[53]:
    Ohio  Nevada  Missouri  Alabama
a      1        2       NaN      NaN
b     NaN      NaN        7        8
c      3        4       9       10
d     NaN      NaN       11       12
e      5        6      13       14
```

In part for legacy reasons (much earlier versions of pandas), DataFrame's `join` method performs a left join on the join keys. It also supports joining the index of the passed DataFrame on one of the columns of the calling DataFrame:

```
In [54]: left1.join(right1, on='key')
Out[54]:
   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
```

Lastly, for simple index-on-index merges, you can pass a list of DataFrames to `join` as an alternative to using the more general `concat` function described below:

```
In [55]: another = DataFrame([[7., 8.], [9., 10.], [11., 12.], [16., 17.]],
....:                      index=['a', 'c', 'e', 'f'], columns=['New York', 'Oregon'])

In [56]: left2.join([right2, another])
Out[56]:
    Ohio  Nevada  Missouri  Alabama  New York  Oregon
a      1        2       NaN      NaN        7        8
c      3        4       9       10        9       10
e      5        6      13       14       11       12

In [57]: left2.join([right2, another], how='outer')
Out[57]:
    Ohio  Nevada  Missouri  Alabama  New York  Oregon
a      1        2       NaN      NaN        7        8
b     NaN      NaN        7        8       NaN      NaN
c      3        4       9       10        9       10
d     NaN      NaN       11       12       NaN      NaN
e      5        6      13       14       11       12
f     NaN      NaN      NaN      NaN       16       17
```

## Concatenating Along an Axis

Another kind of data combination operation is alternatively referred to as concatenation, binding, or stacking. NumPy has a `concatenate` function for doing this with raw NumPy arrays:

```
In [58]: arr = np.arange(12).reshape((3, 4))
```

```
In [59]: arr
```

```
Out[59]:
```

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

```
In [60]: np.concatenate([arr, arr], axis=1)
```

```
Out[60]:
```

```
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]])
```

In the context of pandas objects such as Series and DataFrame, having labeled axes enable you to further generalize array concatenation. In particular, you have a number of additional things to think about:

- If the objects are indexed differently on the other axes, should the collection of axes be unioned or intersected?
- Do the groups need to be identifiable in the resulting object?
- Does the concatenation axis matter at all?

The `concat` function in pandas provides a consistent way to address each of these concerns. I'll give a number of examples to illustrate how it works. Suppose we have three Series with no index overlap:

```
In [61]: s1 = Series([0, 1], index=['a', 'b'])
```

```
In [62]: s2 = Series([2, 3, 4], index=['c', 'd', 'e'])
```

```
In [63]: s3 = Series([5, 6], index=['f', 'g'])
```

Calling `concat` with these object in a list glues together the values and indexes:

```
In [64]: pd.concat([s1, s2, s3])
```

```
Out[64]:
```

```
a    0
b    1
c    2
d    3
e    4
f    5
g    6
```

By default `concat` works along `axis=0`, producing another Series. If you pass `axis=1`, the result will instead be a DataFrame (`axis=1` is the columns):

```
In [65]: pd.concat([s1, s2, s3], axis=1)
Out[65]:
   0   1   2
a  0  NaN  NaN
b  1  NaN  NaN
c NaN  2  NaN
d NaN  3  NaN
e NaN  4  NaN
f NaN  NaN  5
g NaN  NaN  6
```

In this case there is no overlap on the other axis, which as you can see is the sorted union (the '`outer`' join) of the indexes. You can instead intersect them by passing `join='inner'`:

```
In [66]: s4 = pd.concat([s1 * 5, s3])
```

```
In [67]: pd.concat([s1, s4], axis=1)
Out[67]:
   0   1
a  0  0
b  1  5
f NaN  5
g NaN  6
```

```
In [68]: pd.concat([s1, s4], axis=1, join='inner')
Out[68]:
   0   1
a  0  0
b  1  5
```

You can even specify the axes to be used on the other axes with `join_axes`:

```
In [69]: pd.concat([s1, s4], axis=1, join_axes=[['a', 'c', 'b', 'e']])
Out[69]:
   0   1
a  0  0
c NaN  NaN
b  1  5
e NaN  NaN
```

One issue is that the concatenated pieces are not identifiable in the result. Suppose instead you wanted to create a hierarchical index on the concatenation axis. To do this, use the `keys` argument:

```
In [70]: result = pd.concat([s1, s1, s3], keys=['one', 'two', 'three'])
```

```
In [71]: result
Out[71]:
one    a    0
      b    1
two    a    0
      b    1
three   f    5
        g    6

# Much more on the unstack function later
In [72]: result.unstack()
Out[72]:
```

```

      a   b   f   g
one    0   1  NaN  NaN
two    0   1  NaN  NaN
three  NaN  NaN   5   6

```

In the case of combining Series along `axis=1`, the keys become the DataFrame column headers:

```

In [73]: pd.concat([s1, s2, s3], axis=1, keys=['one', 'two', 'three'])
Out[73]:
      one  two  three
a      0   NaN  NaN
b      1   NaN  NaN
c     NaN   2   NaN
d     NaN   3   NaN
e     NaN   4   NaN
f     NaN  NaN   5
g     NaN  NaN   6

```

The same logic extends to DataFrame objects:

```

In [74]: df1 = DataFrame(np.arange(6).reshape(3, 2), index=['a', 'b', 'c'],
....:                  columns=['one', 'two'])

In [75]: df2 = DataFrame(5 + np.arange(4).reshape(2, 2), index=['a', 'c'],
....:                  columns=['three', 'four'])

In [76]: pd.concat([df1, df2], axis=1, keys=['level1', 'level2'])
Out[76]:
      level1      level2
      one  two  three  four
a      0   1      5   6
b      2   3     NaN  NaN
c      4   5      7   8

```

If you pass a dict of objects instead of a list, the dict's keys will be used for the `keys` option:

```

In [77]: pd.concat({'level1': df1, 'level2': df2}, axis=1)
Out[77]:
      level1      level2
      one  two  three  four
a      0   1      5   6
b      2   3     NaN  NaN
c      4   5      7   8

```

There are a couple of additional arguments governing how the hierarchical index is created (see [Table 7-2](#)):

```

In [78]: pd.concat([df1, df2], axis=1, keys=['level1', 'level2'],
....:                 names=['upper', 'lower'])
Out[78]:
      upper  level1      level2
      lower   one  two  three  four
a          0   1      5   6
b          2   3     NaN  NaN
c          4   5      7   8

```

A last consideration concerns DataFrames in which the row index is not meaningful in the context of the analysis:

```
In [79]: df1 = DataFrame(np.random.randn(3, 4), columns=['a', 'b', 'c', 'd'])
```

```
In [80]: df2 = DataFrame(np.random.randn(2, 3), columns=['b', 'd', 'a'])
```

```
In [81]: df1
```

```
Out[81]:
```

	a	b	c	d
0	-0.204708	0.478943	-0.519439	-0.555730
1	1.965781	1.393406	0.092908	0.281746
2	0.769023	1.246435	1.007189	-1.296221

```
In [82]: df2
```

```
Out[82]:
```

	b	d	a
0	0.274992	0.228913	1.352917
1	0.886429	-2.001637	-0.371843

In this case, you can pass `ignore_index=True`:

```
In [83]: pd.concat([df1, df2], ignore_index=True)
```

```
Out[83]:
```

	a	b	c	d
0	-0.204708	0.478943	-0.519439	-0.555730
1	1.965781	1.393406	0.092908	0.281746
2	0.769023	1.246435	1.007189	-1.296221
3	1.352917	0.274992	NaN	0.228913
4	-0.371843	0.886429	NaN	-2.001637

Table 7-2. concat function arguments

Argument	Description
objs	List or dict of pandas objects to be concatenated. The only required argument
axis	Axis to concatenate along; defaults to 0
join	One of 'inner', 'outer', defaulting to 'outer'; whether to intersection (inner) or union (outer) together indexes along the other axes
join_axes	Specific indexes to use for the other n-1 axes instead of performing union/intersection logic
keys	Values to associate with objects being concatenated, forming a hierarchical index along the concatenation axis. Can either be a list or array of arbitrary values, an array of tuples, or a list of arrays (if multiple level arrays passed in levels)
levels	Specific indexes to use as hierarchical index level or levels if keys passed
names	Names for created hierarchical levels if keys and / or levels passed
verify_integrity	Check new axis in concatenated object for duplicates and raise exception if so. By default (False) allows duplicates
ignore_index	Do not preserve indexes along concatenation axis, instead producing a new range(total_length) index

## Combining Data with Overlap

Another data combination situation can't be expressed as either a merge or concatenation operation. You may have two datasets whose indexes overlap in full or part. As a motivating example, consider NumPy's `where` function, which expressed a vectorized if-else:

```
In [84]: a = Series([np.nan, 2.5, np.nan, 3.5, 4.5, np.nan],
....:                  index=['f', 'e', 'd', 'c', 'b', 'a'])

In [85]: b = Series(np.arange(len(a), dtype=np.float64),
....:                  index=['f', 'e', 'd', 'c', 'b', 'a'])

In [86]: b[-1] = np.nan

In [87]: a      In [88]: b      In [89]: np.where(pd.isnull(a), b, a)
Out[87]:      Out[88]:      Out[89]:
f    NaN      f    0      f    0.0
e    2.5      e    1      e    2.5
d    NaN      d    2      d    2.0
c    3.5      c    3      c    3.5
b    4.5      b    4      b    4.5
a    NaN      a    NaN    a    NaN
```

Series has a `combine_first` method, which performs the equivalent of this operation plus data alignment:

```
In [90]: b[:-2].combine_first(a[2:])
Out[90]:
a    NaN
b    4.5
c    3.0
d    2.0
e    1.0
f    0.0
```

With DataFrames, `combine_first` naturally does the same thing column by column, so you can think of it as “patching” missing data in the calling object with data from the object you pass:

```
In [91]: df1 = DataFrame({'a': [1., np.nan, 5., np.nan],
....:                      'b': [np.nan, 2., np.nan, 6.],
....:                      'c': range(2, 18, 4)})

In [92]: df2 = DataFrame({'a': [5., 4., np.nan, 3., 7.],
....:                      'b': [np.nan, 3., 4., 6., 8.]})

In [93]: df1.combine_first(df2)
Out[93]:
   a   b   c
0  1  NaN  2
1  4   2   6
2  5   4  10
3  3   6  14
4  7   8  NaN
```

## Reshaping and Pivoting

There are a number of fundamental operations for rearranging tabular data. These are alternately referred to as *reshape* or *pivot* operations.

## Reshaping with Hierarchical Indexing

Hierarchical indexing provides a consistent way to rearrange data in a DataFrame. There are two primary actions:

- `stack`: this “rotates” or pivots from the columns in the data to the rows
- `unstack`: this pivots from the rows into the columns

I’ll illustrate these operations through a series of examples. Consider a small DataFrame with string arrays as row and column indexes:

```
In [94]: data = DataFrame(np.arange(6).reshape((2, 3)),  
....:                      index=pd.Index(['Ohio', 'Colorado'], name='state'),  
....:                      columns=pd.Index(['one', 'two', 'three'], name='number'))  
  
In [95]: data  
Out[95]:  
number   one   two   three  
state  
Ohio      0     1     2  
Colorado  3     4     5
```

Using the `stack` method on this data pivots the columns into the rows, producing a Series:

```
In [96]: result = data.stack()  
  
In [97]: result  
Out[97]:  
state   number  
Ohio    one      0  
        two      1  
        three    2  
Colorado one      3  
        two      4  
        three    5
```

From a hierarchically-indexed Series, you can rearrange the data back into a DataFrame with `unstack`:

```
In [98]: result.unstack()  
Out[98]:  
number   one   two   three  
state  
Ohio      0     1     2  
Colorado  3     4     5
```

By default the innermost level is unstacked (same with `stack`). You can unstack a different level by passing a level number or name:

```
In [99]: result.unstack(0)          In [100]: result.unstack('state')  
Out[99]:  
state  Ohio  Colorado  
number  
one      0       3  
                                Out[100]:  
                                state  Ohio  Colorado  
                                number  
                                one      0       3
```

```

two      1      4      two      1      4
three    2      5      three    2      5

```

Unstacking might introduce missing data if all of the values in the level aren't found in each of the subgroups:

```

In [101]: s1 = Series([0, 1, 2, 3], index=['a', 'b', 'c', 'd'])

In [102]: s2 = Series([4, 5, 6], index=['c', 'd', 'e'])

In [103]: data2 = pd.concat([s1, s2], keys=['one', 'two'])

In [104]: data2.unstack()
Out[104]:
   a   b   c   d   e
one  0   1   2   3  NaN
two  NaN NaN  4   5   6

```

Stacking filters out missing data by default, so the operation is easily invertible:

```

In [105]: data2.unstack().stack()      In [106]: data2.unstack().stack(dropna=False)
Out[105]:                                Out[106]:
one  a      0                         one  a      0
      b      1                         b      1
      c      2                         c      2
      d      3                         d      3
two  c      4                         e      NaN
      d      5                         two  a      NaN
      e      6                         b      NaN
                                         c      4
                                         d      5
                                         e      6

```

When unstacking in a DataFrame, the level unstacked becomes the lowest level in the result:

```

In [107]: df = DataFrame({'left': result, 'right': result + 5},
.....:                      columns=pd.Index(['left', 'right'], name='side'))

In [108]: df
Out[108]:
   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 [109]: df.unstack('state')          In [110]: df.unstack('state').stack('side')
Out[109]:                                Out[110]:
   side      left      right
state  Ohio  Colorado  Ohio  Colorado
   number
   one      0        3      5      8
   two      1        4      6      9

```

	state	Ohio	Colorado	state	Ohio	Colorado
number				number	side	
one	left	0	3	one	left	0
two	right	5	8	two	right	5

three	2	5	7	10		right	6	9
			three		left	2	5	

				right	7	10	
--	--	--	--	-------	---	----	--

## Pivoting “long” to “wide” Format

A common way to store multiple time series in databases and CSV is in so-called *long* or *stacked* format:

```
In [116]: ldata[:10]
Out[116]:
      date    item    value
0 1959-03-31 00:00:00  realgdp  2710.349
1 1959-03-31 00:00:00     infl    0.000
2 1959-03-31 00:00:00    unemp   5.800
3 1959-06-30 00:00:00  realgdp  2778.801
4 1959-06-30 00:00:00     infl   2.340
5 1959-06-30 00:00:00    unemp   5.100
6 1959-09-30 00:00:00  realgdp  2775.488
7 1959-09-30 00:00:00     infl   2.740
8 1959-09-30 00:00:00    unemp   5.300
9 1959-12-31 00:00:00  realgdp  2785.204
```

Data is frequently stored this way in relational databases like MySQL as a fixed schema (column names and data types) allows the number of distinct values in the `item` column to increase or decrease as data is added or deleted in the table. In the above example `date` and `item` would usually be the primary keys (in relational database parlance), offering both relational integrity and easier joins and programmatic queries in many cases. The downside, of course, is that the data may not be easy to work with in long format; you might prefer to have a DataFrame containing one column per distinct `item` value indexed by timestamps in the `date` column. DataFrame’s `pivot` method performs exactly this transformation:

```
In [117]: pivoted = ldata.pivot('date', 'item', 'value')

In [118]: pivoted.head()
Out[118]:
      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
```

The first two values passed are the columns to be used as the row and column index, and finally an optional value column to fill the DataFrame. Suppose you had two value columns that you wanted to reshape simultaneously:

```
In [119]: ldata['value2'] = np.random.randn(len(ldata))

In [120]: ldata[:10]
Out[120]:
```

```

      date      item    value   value2
0 1959-03-31 00:00:00  realgdp  2710.349  1.669025
1 1959-03-31 00:00:00     infl     0.000 -0.438570
2 1959-03-31 00:00:00    unemp     5.800 -0.539741
3 1959-06-30 00:00:00  realgdp  2778.801  0.476985
4 1959-06-30 00:00:00     infl     2.340  3.248944
5 1959-06-30 00:00:00    unemp     5.100 -1.021228
6 1959-09-30 00:00:00  realgdp  2775.488 -0.577087
7 1959-09-30 00:00:00     infl     2.740  0.124121
8 1959-09-30 00:00:00    unemp     5.300  0.302614
9 1959-12-31 00:00:00  realgdp  2785.204  0.523772

```

By omitting the last argument, you obtain a DataFrame with hierarchical columns:

```
In [121]: pivoted = ldata.pivot('date', 'item')
```

```
In [122]: pivoted[:5]
```

```
Out[122]:
```

item	value			value2		
	infl	realgdp	unemp	infl	realgdp	unemp
date						
1959-03-31	0.00	2710.349	5.8	-0.438570	1.669025	-0.539741
1959-06-30	2.34	2778.801	5.1	3.248944	0.476985	-1.021228
1959-09-30	2.74	2775.488	5.3	0.124121	-0.577087	0.302614
1959-12-31	0.27	2785.204	5.6	0.000940	0.523772	1.343810
1960-03-31	2.31	2847.699	5.2	-0.831154	-0.713544	-2.370232

```
In [123]: pivoted['value'][:5]
```

```
Out[123]:
```

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

Note that `pivot` is just a shortcut for creating a hierarchical index using `set_index` and reshaping with `unstack`:

```
In [124]: unstacked = ldata.set_index(['date', 'item']).unstack('item')
```

```
In [125]: unstacked[:7]
```

```
Out[125]:
```

item	value			value2		
	infl	realgdp	unemp	infl	realgdp	unemp
date						
1959-03-31	0.00	2710.349	5.8	-0.438570	1.669025	-0.539741
1959-06-30	2.34	2778.801	5.1	3.248944	0.476985	-1.021228
1959-09-30	2.74	2775.488	5.3	0.124121	-0.577087	0.302614
1959-12-31	0.27	2785.204	5.6	0.000940	0.523772	1.343810
1960-03-31	2.31	2847.699	5.2	-0.831154	-0.713544	-2.370232
1960-06-30	0.14	2834.390	5.2	-0.860757	-1.860761	0.560145
1960-09-30	2.70	2839.022	5.6	0.119827	-1.265934	-1.063512

# Data Transformation

So far in this chapter we've been concerned with rearranging data. Filtering, cleaning, and other transformations are another class of important operations.

## Removing Duplicates

Duplicate rows may be found in a DataFrame for any number of reasons. Here is an example:

```
In [126]: data = DataFrame({'k1': ['one'] * 3 + ['two'] * 4,
.....:                  'k2': [1, 1, 2, 3, 3, 4, 4]})

In [127]: data
Out[127]:
   k1  k2
0  one  1
1  one  1
2  one  2
3  two  3
4  two  3
5  two  4
6  two  4
```

The DataFrame method `duplicated` returns a boolean Series indicating whether each row is a duplicate or not:

```
In [128]: data.duplicated()
Out[128]:
0    False
1     True
2    False
3    False
4     True
5    False
6     True
```

Relatedly, `drop_duplicates` returns a DataFrame where the `duplicated` array is `True`:

```
In [129]: data.drop_duplicates()
Out[129]:
   k1  k2
0  one  1
2  one  2
3  two  3
5  two  4
```

Both of these methods by default consider all of the columns; alternatively you can specify any subset of them to detect duplicates. Suppose we had an additional column of values and wanted to filter duplicates only based on the `'k1'` column:

```
In [130]: data['v1'] = range(7)

In [131]: data.drop_duplicates(['k1'])
```

```
Out[131]:  
      k1  k2  v1  
0   one    1    0  
3   two    3    3
```

duplicated and drop\_duplicates by default keep the first observed value combination. Passing take\_last=True will return the last one:

```
In [132]: data.drop_duplicates(['k1', 'k2'], take_last=True)  
Out[132]:  
      k1  k2  v1  
1   one    1    1  
2   one    2    2  
4   two    3    4  
6   two    4    6
```

## Transforming Data Using a Function or Mapping

For many data sets, you may wish to perform some transformation based on the values in an array, Series, or column in a DataFrame. Consider the following hypothetical data collected about some kinds of meat:

```
In [133]: data = DataFrame({'food': ['bacon', 'pulled pork', 'bacon', 'Pastrami',  
.....:                               'corned beef', 'Bacon', 'pastrami', 'honey ham',  
.....:                               'nova lox'],  
.....:                               'ounces': [4, 3, 12, 6, 7.5, 8, 3, 5, 6]})  
  
In [134]: data  
Out[134]:  
      food  ounces  
0     bacon    4.0  
1  pulled pork    3.0  
2     bacon   12.0  
3    Pastrami    6.0  
4  corned beef    7.5  
5      Bacon    8.0  
6    pastrami    3.0  
7  honey ham    5.0  
8   nova lox    6.0
```

Suppose you wanted to add a column indicating the type of animal that each food came from. Let's write down a mapping of each distinct meat type to the kind of animal:

```
meat_to_animal = {  
    'bacon': 'pig',  
    'pulled pork': 'pig',  
    'pastrami': 'cow',  
    'corned beef': 'cow',  
    'honey ham': 'pig',  
    'nova lox': 'salmon'  
}
```

The `map` method on a Series accepts a function or dict-like object containing a mapping, but here we have a small problem in that some of the meats above are capitalized and others are not. Thus, we also need to convert each value to lower case:

```
In [136]: data['animal'] = data['food'].map(str.lower).map(meat_to_animal)
```

```
In [137]: data
```

```
Out[137]:
```

	food	ounces	animal
0	bacon	4.0	pig
1	pulled pork	3.0	pig
2	bacon	12.0	pig
3	Pastrami	6.0	cow
4	corned beef	7.5	cow
5	Bacon	8.0	pig
6	pastrami	3.0	cow
7	honey ham	5.0	pig
8	nova lox	6.0	salmon

We could also have passed a function that does all the work:

```
In [138]: data['food'].map(lambda x: meat_to_animal[x.lower()])
```

```
Out[138]:
```

0	pig
1	pig
2	pig
3	cow
4	cow
5	pig
6	cow
7	pig
8	salmon

```
Name: food
```

Using `map` is a convenient way to perform element-wise transformations and other data cleaning-related operations.

## Replacing Values

Filling in missing data with the `fillna` method can be thought of as a special case of more general value replacement. While `map`, as you've seen above, can be used to modify a subset of values in an object, `replace` provides a simpler and more flexible way to do so. Let's consider this Series:

```
In [139]: data = Series([1., -999., 2., -999., -1000., 3.])
```

```
In [140]: data
```

```
Out[140]:
```

0	1
1	-999
2	2
3	-999
4	-1000
5	3

The `-999` values might be sentinel values for missing data. To replace these with NA values that pandas understands, we can use `replace`, producing a new Series:

```
In [141]: data.replace(-999, np.nan)
Out[141]:
0      1
1    NaN
2      2
3    NaN
4   -1000
5      3
```

If you want to replace multiple values at once, you instead pass a list then the substitute value:

```
In [142]: data.replace([-999, -1000], np.nan)
Out[142]:
0      1
1    NaN
2      2
3    NaN
4    NaN
5      3
```

To use a different replacement for each value, pass a list of substitutes:

```
In [143]: data.replace([-999, -1000], [np.nan, 0])
Out[143]:
0      1
1    NaN
2      2
3    NaN
4      0
5      3
```

The argument passed can also be a dict:

```
In [144]: data.replace({-999: np.nan, -1000: 0})
Out[144]:
0      1
1    NaN
2      2
3    NaN
4      0
5      3
```

## Renaming Axis Indexes

Like values in a Series, axis labels can be similarly transformed by a function or mapping of some form to produce new, differently labeled objects. The axes can also be modified in place without creating a new data structure. Here's a simple example:

```
In [145]: data = DataFrame(np.arange(12).reshape((3, 4)),
.....:                   index=['Ohio', 'Colorado', 'New York'],
.....:                   columns=['one', 'two', 'three', 'four'])
```

Like a Series, the axis indexes have a `map` method:

```
In [146]: data.index.map(str.upper)
Out[146]: array([OHIO, COLORADO, NEW YORK], dtype=object)
```

You can assign to `index`, modifying the DataFrame in place:

```
In [147]: data.index = data.index.map(str.upper)
```

```
In [148]: data
Out[148]:
   one  two  three  four
OHIO    0    1      2    3
COLORADO 4    5      6    7
NEW YORK 8    9     10    11
```

If you want to create a transformed version of a data set without modifying the original, a useful method is `rename`:

```
In [149]: data.rename(index=str.title, columns=str.upper)
Out[149]:
   ONE  TWO  THREE  FOUR
Ohio    0    1      2    3
Colorado 4    5      6    7
New York 8    9     10    11
```

Notably, `rename` can be used in conjunction with a dict-like object providing new values for a subset of the axis labels:

```
In [150]: data.rename(index={'OHIO': 'INDIANA'},
                     ....:                 columns={'three': 'peekaboo'})
Out[150]:
   one  two  peekaboo  four
INDIANA  0    1        2    3
COLORADO 4    5        6    7
NEW YORK 8    9       10    11
```

`rename` saves having to copy the DataFrame manually and assign to its `index` and `columns` attributes. Should you wish to modify a data set in place, pass `inplace=True`:

```
# Always returns a reference to a DataFrame
In [151]: _ = data.rename(index={'OHIO': 'INDIANA'}, inplace=True)

In [152]: data
Out[152]:
   one  two  three  four
INDIANA  0    1      2    3
COLORADO 4    5      6    7
NEW YORK 8    9     10    11
```

## Discretization and Binning

Continuous data is often discretized or otherwise separated into “bins” for analysis. Suppose you have data about a group of people in a study, and you want to group them into discrete age buckets:

```
In [153]: ages = [20, 22, 25, 27, 21, 23, 37, 31, 61, 45, 41, 32]
```

Let’s divide these into bins of 18 to 25, 26 to 35, 35 to 60, and finally 60 and older. To do so, you have to use `cut`, a function in pandas:

```
In [154]: bins = [18, 25, 35, 60, 100]
```

```
In [155]: cats = pd.cut(ages, bins)
```

```
In [156]: cats
```

```
Out[156]:
```

```
Categorical:
```

```
array([(18, 25], (18, 25], (18, 25], (25, 35], (18, 25], (18, 25],
       (35, 60], (25, 35], (60, 100], (35, 60], (35, 60], (25, 35]], dtype=object)
Levels (4): Index([(18, 25], (25, 35], (35, 60], (60, 100]], dtype=object)
```

The object pandas returns is a special `Categorical` object. You can treat it like an array of strings indicating the bin name; internally it contains a `levels` array indicating the distinct category names along with a labeling for the `ages` data in the `labels` attribute:

```
In [157]: cats.labels
```

```
Out[157]: array([0, 0, 0, 1, 0, 0, 2, 1, 3, 2, 2, 1])
```

```
In [158]: cats.levels
```

```
Out[158]: Index([(18, 25], (25, 35], (35, 60], (60, 100]], dtype=object)
```

```
In [159]: pd.value_counts(cats)
```

```
Out[159]:
```

(18, 25]	5
(35, 60]	3
(25, 35]	3
(60, 100]	1

Consistent with mathematical notation for intervals, a parenthesis means that the side is *open* while the square bracket means it is *closed* (inclusive). Which side is closed can be changed by passing `right=False`:

```
In [160]: pd.cut(ages, [18, 26, 36, 61, 100], right=False)
```

```
Out[160]:
```

```
Categorical:
```

```
array([(18, 26), [18, 26), [18, 26), [26, 36), [18, 26), [18, 26),
       [36, 61), [26, 36), [61, 100], [36, 61), [36, 61), [26, 36]], dtype=object)
Levels (4): Index([(18, 26), [26, 36), [36, 61), [61, 100]], dtype=object)
```

You can also pass your own bin names by passing a list or array to the `labels` option:

```
In [161]: group_names = ['Youth', 'YoungAdult', 'MiddleAged', 'Senior']
```

```
In [162]: pd.cut(ages, bins, labels=group_names)
```

```
Out[162]:
```

```
Categorical:  
array([Youth, Youth, Youth, YoungAdult, Youth, Youth, MiddleAged,  
      YoungAdult, Senior, MiddleAged, MiddleAged, YoungAdult], dtype=object)  
Levels (4): Index([Youth, YoungAdult, MiddleAged, Senior], dtype=object)
```

If you pass `cut` a integer number of bins instead of explicit bin edges, it will compute equal-length bins based on the minimum and maximum values in the data. Consider the case of some uniformly distributed data chopped into fourths:

```
In [163]: data = np.random.rand(20)  
  
In [164]: pd.cut(data, 4, precision=2)  
Out[164]:  
Categorical:  
array([(0.45, 0.67], (0.23, 0.45], (0.0037, 0.23], (0.45, 0.67],  
      (0.67, 0.9], (0.45, 0.67], (0.67, 0.9], (0.23, 0.45], (0.23, 0.45],  
      (0.67, 0.9], (0.67, 0.9], (0.67, 0.9], (0.23, 0.45], (0.23, 0.45],  
      (0.23, 0.45], (0.67, 0.9], (0.0037, 0.23], (0.0037, 0.23],  
      (0.23, 0.45], (0.23, 0.45]], dtype=object)  
Levels (4): Index([(0.0037, 0.23], (0.23, 0.45], (0.45, 0.67],  
                  (0.67, 0.9]], dtype=object)
```

A closely related function, `qcut`, bins the data based on sample quantiles. Depending on the distribution of the data, using `cut` will not usually result in each bin having the same number of data points. Since `qcut` uses sample quantiles instead, by definition you will obtain roughly equal-size bins:

```
In [165]: data = np.random.randn(1000) # Normally distributed  
  
In [166]: cats = pd.qcut(data, 4) # Cut into quartiles  
  
In [167]: cats  
Out[167]:  
Categorical:  
array([[-0.022, 0.641], [-3.745, -0.635], (0.641, 3.26], ...,  
      (-0.635, -0.022], (0.641, 3.26], (-0.635, -0.022]], dtype=object)  
Levels (4): Index([-3.745, -0.635], (-0.635, -0.022], (-0.022, 0.641],  
                  (0.641, 3.26]], dtype=object)  
  
In [168]: pd.value_counts(cats)  
Out[168]:  
[-3.745, -0.635]    250  
(0.641, 3.26]      250  
(-0.635, -0.022]   250  
(-0.022, 0.641]    250
```

Similar to `cut` you can pass your own quantiles (numbers between 0 and 1, inclusive):

```
In [169]: pd.qcut(data, [0, 0.1, 0.5, 0.9, 1.])  
Out[169]:  
Categorical:  
array([[-0.022, 1.302], (-1.266, -0.022], (-0.022, 1.302], ...,  
      (-1.266, -0.022], (-0.022, 1.302], (-1.266, -0.022]], dtype=object)  
Levels (4): Index([-3.745, -1.266], (-1.266, -0.022], (-0.022, 1.302],  
                  (1.302, 3.26]], dtype=object)
```

We'll return to `cut` and `qcut` later in the chapter on aggregation and group operations, as these discretization functions are especially useful for quantile and group analysis.

## Detecting and Filtering Outliers

Filtering or transforming outliers is largely a matter of applying array operations. Consider a DataFrame with some normally distributed data:

```
In [170]: np.random.seed(12345)
```

```
In [171]: data = DataFrame(np.random.randn(1000, 4))
```

```
In [172]: data.describe()
```

```
Out[172]:
```

	0	1	2	3
count	1000.000000	1000.000000	1000.000000	1000.000000
mean	-0.067684	0.067924	0.025598	-0.002298
std	0.998035	0.992106	1.006835	0.996794
min	-3.428254	-3.548824	-3.184377	-3.745356
25%	-0.774890	-0.591841	-0.641675	-0.644144
50%	-0.116401	0.101143	0.002073	-0.013611
75%	0.616366	0.780282	0.680391	0.654328
max	3.366626	2.653656	3.260383	3.927528

Suppose you wanted to find values in one of the columns exceeding three in magnitude:

```
In [173]: col = data[3]
```

```
In [174]: col[np.abs(col) > 3]
```

```
Out[174]:
```

97	3.927528
305	-3.399312
400	-3.745356
Name:	3

To select all rows having a value exceeding 3 or -3, you can use the `any` method on a boolean DataFrame:

```
In [175]: data[(np.abs(data) > 3).any(1)]
```

```
Out[175]:
```

	0	1	2	3
5	-0.539741	0.476985	3.248944	-1.021228
97	-0.774363	0.552936	0.106061	3.927528
102	-0.655054	-0.565230	3.176873	0.959533
305	-2.315555	0.457246	-0.025907	-3.399312
324	0.050188	1.951312	3.260383	0.963301
400	0.146326	0.508391	-0.196713	-3.745356
499	-0.293333	-0.242459	-3.056990	1.918403
523	-3.428254	-0.296336	-0.439938	-0.867165
586	0.275144	1.179227	-3.184377	1.369891
808	-0.362528	-3.548824	1.553205	-2.186301
900	3.366626	-2.372214	0.851010	1.332846

Values can just as easily be set based on these criteria. Here is code to cap values outside the interval -3 to 3:

```
In [176]: data[np.abs(data) > 3] = np.sign(data) * 3
```

```
In [177]: data.describe()
```

```
Out[177]:
```

	0	1	2	3
count	1000.000000	1000.000000	1000.000000	1000.000000
mean	-0.067623	0.068473	0.025153	-0.002081
std	0.995485	0.990253	1.003977	0.989736
min	-3.000000	-3.000000	-3.000000	-3.000000
25%	-0.774890	-0.591841	-0.641675	-0.644144
50%	-0.116401	0.101143	0.002073	-0.013611
75%	0.616366	0.780282	0.680391	0.654328
max	3.000000	2.653656	3.000000	3.000000

The ufunc `np.sign` returns an array of 1 and -1 depending on the sign of the values.

## Permutation and Random Sampling

Permuting (randomly reordering) a Series or the rows in a DataFrame is easy to do using the `numpy.random.permutation` function. Calling `permutation` with the length of the axis you want to permute produces an array of integers indicating the new ordering:

```
In [178]: df = DataFrame(np.arange(5 * 4).reshape(5, 4))
```

```
In [179]: sampler = np.random.permutation(5)
```

```
In [180]: sampler
```

```
Out[180]: array([1, 0, 2, 3, 4])
```

That array can then be used in `ix`-based indexing or the `take` function:

```
In [181]: df
```

```
Out[181]:
```

0	1	2	3
0	0	1	2
1	4	5	6
2	8	9	10
3	12	13	14
4	16	17	18

```
In [182]: df.take(sampler)
```

```
Out[182]:
```

0	1	2	3
1	4	5	6
0	0	1	2
2	8	9	10
3	12	13	14
4	16	17	18

To select a random subset without replacement, one way is to slice off the first  $k$  elements of the array returned by `permutation`, where  $k$  is the desired subset size. There are much more efficient sampling-without-replacement algorithms, but this is an easy strategy that uses readily available tools:

```
In [183]: df.take(np.random.permutation(len(df))[:3])
```

```
Out[183]:
```

0	1	2	3
1	4	5	6
3	12	13	14
4	16	17	18

To generate a sample *with* replacement, the fastest way is to use `np.random.randint` to draw random integers:

```
In [184]: bag = np.array([5, 7, -1, 6, 4])  
In [185]: sampler = np.random.randint(0, len(bag), size=10)  
In [186]: sampler  
Out[186]: array([4, 4, 2, 2, 2, 0, 3, 0, 4, 1])  
In [187]: draws = bag.take(sampler)  
In [188]: draws  
Out[188]: array([ 4,  4, -1, -1,  5,  6,  5,  4,  7])
```

## Computing Indicator/Dummy Variables

Another type of transformation for statistical modeling or machine learning applications is converting a categorical variable into a “dummy” or “indicator” matrix. If a column in a DataFrame has k distinct values, you would derive a matrix or DataFrame containing k columns containing all 1’s and 0’s. pandas has a `get_dummies` function for doing this, though devising one yourself is not difficult. Let’s return to an earlier example DataFrame:

```
In [189]: df = DataFrame({'key': ['b', 'b', 'a', 'c', 'a', 'b'],  
.....: 'data1': range(6)})  
In [190]: pd.get_dummies(df['key'])  
Out[190]:  
   a   b   c  
0  0   1   0  
1  0   1   0  
2  1   0   0  
3  0   0   1  
4  1   0   0  
5  0   1   0
```

In some cases, you may want to add a prefix to the columns in the indicator DataFrame, which can then be merged with the other data. `get_dummies` has a `prefix` argument for doing just this:

```
In [191]: dummies = pd.get_dummies(df['key'], prefix='key')  
In [192]: df_with_dummy = df[['data1']].join(dummies)  
In [193]: df_with_dummy  
Out[193]:  
   data1  key_a  key_b  key_c  
0      0      0      1      0  
1      1      0      1      0  
2      2      1      0      0  
3      3      0      0      1  
4      4      1      0      0  
5      5      0      1      0
```

If a row in a DataFrame belongs to multiple categories, things are a bit more complicated. Let's return to the MovieLens 1M dataset from earlier in the book:

```
In [194]: mnames = ['movie_id', 'title', 'genres']

In [195]: movies = pd.read_table('ch07/movies.dat', sep='::', header=None,
.....:             names=mnames)

In [196]: movies[:10]
Out[196]:
   movie_id          title           genres
0        1      Toy Story (1995) Animation|Children's|Comedy
1        2       Jumanji (1995) Adventure|Children's|Fantasy
2        3  Grumpier Old Men (1995)            Comedy|Romance
3        4    Waiting to Exhale (1995)            Comedy|Drama
4        5 Father of the Bride Part II (1995)            Comedy
5        6              Heat (1995) Action|Crime|Thriller
6        7            Sabrina (1995)            Comedy|Romance
7        8        Tom and Huck (1995) Adventure|Children's
8        9     Sudden Death (1995)            Action
9       10        GoldenEye (1995) Action|Adventure|Thriller
```

Adding indicator variables for each genre requires a little bit of wrangling. First, we extract the list of unique genres in the dataset (using a nice `set.union` trick):

```
In [197]: genre_iter = (set(x.split('|')) for x in movies.genres)

In [198]: genres = sorted(set.union(*genre_iter))
```

Now, one way to construct the indicator DataFrame is to start with a DataFrame of all zeros:

```
In [199]: dummies = DataFrame(np.zeros((len(movies), len(genres))), columns=genres)
```

Now, iterate through each movie and set entries in each row of `dummies` to 1:

```
In [200]: for i, gen in enumerate(movies.genres):
.....:     dummies.ix[i, gen.split('|')] = 1
```

Then, as above, you can combine this with `movies`:

```
In [201]: movies_windic = movies.join(dummies.add_prefix('Genre_'))

In [202]: movies_windic.ix[0]
Out[202]:
   movie_id          title           1
   genres           Animation|Children's|Comedy
   Genre_Action     0
   Genre_Adventure 0
   Genre_Animation  1
   Genre_Children's 1
   Genre_Comedy     1
   Genre_Crime      0
   Genre_Documentary 0
   Genre_Drama      0
   Genre_Fantasy    0
```

```
Genre_Film-Noir          0
Genre_Horror              0
Genre_Musical             0
Genre_Mystery             0
Genre_Romance             0
Genre_Sci-Fi              0
Genre_Thriller            0
Genre_War                 0
Genre_Western              Name: 0
```



For much larger data, this method of constructing indicator variables with multiple membership is not especially speedy. A lower-level function leveraging the internals of the DataFrame could certainly be written.

A useful recipe for statistical applications is to combine `get_dummies` with a discretization function like `cut`:

```
In [204]: values = np.random.rand(10)
```

```
In [205]: values
```

```
Out[205]:
```

```
array([ 0.9296,  0.3164,  0.1839,  0.2046,  0.5677,  0.5955,  0.9645,
       0.6532,  0.7489,  0.6536])
```

```
In [206]: bins = [0, 0.2, 0.4, 0.6, 0.8, 1]
```

```
In [207]: pd.get_dummies(pd.cut(values, bins))
```

```
Out[207]:
```

	(0, 0.2]	(0.2, 0.4]	(0.4, 0.6]	(0.6, 0.8]	(0.8, 1]
0	0	0	0	0	1
1	0	1	0	0	0
2	1	0	0	0	0
3	0	1	0	0	0
4	0	0	1	0	0
5	0	0	1	0	0
6	0	0	0	0	1
7	0	0	0	1	0
8	0	0	0	1	0
9	0	0	0	1	0

## String Manipulation

Python has long been a popular data munging language in part due to its ease-of-use for string and text processing. Most text operations are made simple with the string object's built-in methods. For more complex pattern matching and text manipulations, regular expressions may be needed. pandas adds to the mix by enabling you to apply string and regular expressions concisely on whole arrays of data, additionally handling the annoyance of missing data.

## String Object Methods

In many string munging and scripting applications, built-in string methods are sufficient. As an example, a comma-separated string can be broken into pieces with `split`:

```
In [208]: val = 'a,b, guido'
```

```
In [209]: val.split(',')
Out[209]: ['a', 'b', ' guido']
```

`split` is often combined with `strip` to trim whitespace (including newlines):

```
In [210]: pieces = [x.strip() for x in val.split(',')]
```

```
In [211]: pieces
Out[211]: ['a', 'b', 'guido']
```

These substrings could be concatenated together with a two-colon delimiter using addition:

```
In [212]: first, second, third = pieces
```

```
In [213]: first + '::' + second + '::' + third
Out[213]: 'a::b::guido'
```

But, this isn't a practical generic method. A faster and more Pythonic way is to pass a list or tuple to the `join` method on the string `::`:

```
In [214]: '::'.join(pieces)
Out[214]: 'a::b::guido'
```

Other methods are concerned with locating substrings. Using Python's `in` keyword is the best way to detect a substring, though `index` and `find` can also be used:

```
In [215]: 'guido' in val
Out[215]: True
```

```
In [216]: val.index(',')
Out[216]: 1
```

```
In [217]: val.find(':')
Out[217]: -1
```

Note the difference between `find` and `index` is that `index` raises an exception if the string isn't found (versus returning -1):

```
In [218]: val.index(':')
-----
ValueError                                Traceback (most recent call last)
<ipython-input-218-280f8b2856ce> in <module>()
      1 val.index(':')
ValueError: substring not found
```

Relatedly, `count` returns the number of occurrences of a particular substring:

```
In [219]: val.count(',')
Out[219]: 2
```

`replace` will substitute occurrences of one pattern for another. This is commonly used to delete patterns, too, by passing an empty string:

```
In [220]: val.replace(',', '::')
Out[220]: 'a::b:: guido'
```

```
In [221]: val.replace(',', '')
Out[221]: 'ab guido'
```

Regular expressions can also be used with many of these operations as you'll see below.

Table 7-3. Python built-in string methods

Argument	Description
count	Return the number of non-overlapping occurrences of substring in the string.
endswith, startswith	Returns True if string ends with suffix (starts with prefix).
join	Use string as delimiter for concatenating a sequence of other strings.
index	Return position of first character in substring if found in the string. Raises ValueError if not found.
find	Return position of first character of <i>first</i> occurrence of substring in the string. Like index, but returns -1 if not found.
rfind	Return position of first character of <i>last</i> occurrence of substring in the string. Returns -1 if not found.
replace	Replace occurrences of string with another string.
strip, rstrip, lstrip	Trim whitespace, including newlines; equivalent to x.strip() (and rstrip, lstrip, respectively) for each element.
split	Break string into list of substrings using passed delimiter.
lower, upper	Convert alphabet characters to lowercase or uppercase, respectively.
ljust, rjust	Left justify or right justify, respectively. Pad opposite side of string with spaces (or some other fill character) to return a string with a minimum width.

## Regular expressions

Regular expressions provide a flexible way to search or match string patterns in text. A single expression, commonly called a *regex*, is a string formed according to the regular expression language. Python's built-in `re` module is responsible for applying regular expressions to strings; I'll give a number of examples of its use here.



The art of writing regular expressions could be a chapter of its own and thus is outside the book's scope. There are many excellent tutorials and references on the internet, such as Zed Shaw's *Learn Regex The Hard Way* (<http://regex.learncodethehardway.org/book/>).

The `re` module functions fall into three categories: pattern matching, substitution, and splitting. Naturally these are all related; a regex describes a pattern to locate in the text, which can then be used for many purposes. Let's look at a simple example: suppose I wanted to split a string with a variable number of whitespace characters (tabs, spaces, and newlines). The regex describing one or more whitespace characters is `\s+`:

```
In [222]: import re  
  
In [223]: text = "foo    bar\t baz  \tqux"  
  
In [224]: re.split('\s+', text)  
Out[224]: ['foo', 'bar', 'baz', 'qux']
```

When you call `re.split(' \s+', text)`, the regular expression is first *compiled*, then its `split` method is called on the passed text. You can compile the regex yourself with `re.compile`, forming a reusable regex object:

```
In [225]: regex = re.compile(' \s+')  
  
In [226]: regex.split(text)  
Out[226]: ['foo', 'bar', 'baz', 'qux']
```

If, instead, you wanted to get a list of all patterns matching the regex, you can use the `findall` method:

```
In [227]: regex.findall(text)  
Out[227]: [' ', '\t ', '\t']
```



To avoid unwanted escaping with `\` in a regular expression, use *raw* string literals like `r'C:\x'` instead of the equivalent `'C:\\x'`.

Creating a regex object with `re.compile` is highly recommended if you intend to apply the same expression to many strings; doing so will save CPU cycles.

`match` and `search` are closely related to `findall`. While `findall` returns all matches in a string, `search` returns only the first match. More rigidly, `match` *only* matches at the beginning of the string. As a less trivial example, let's consider a block of text and a regular expression capable of identifying most email addresses:

```
text = """Dave dave@google.com  
Steve steve@gmail.com  
Rob rob@gmail.com  
Ryan ryan@yahoo.com  
"""\n\npattern = r'[A-Z0-9._%+-]+@[A-Z0-9.-]+\.[A-Z]{2,4}'\n\n# re.IGNORECASE makes the regex case-insensitive\nregex = re.compile(pattern, flags=re.IGNORECASE)
```

Using `findall` on the text produces a list of the e-mail addresses:

```
In [229]: regex.findall(text)  
Out[229]: ['dave@google.com', 'steve@gmail.com', 'rob@gmail.com', 'ryan@yahoo.com']
```

`search` returns a special match object for the first email address in the text. For the above regex, the match object can only tell us the start and end position of the pattern in the string:

```
In [230]: m = regex.search(text)

In [231]: m
Out[231]: <_sre.SRE_Match at 0x10a05de00>

In [232]: text[m.start():m.end()]
Out[232]: 'dave@google.com'
```

`regex.match` returns `None`, as it only will match if the pattern occurs at the start of the string:

```
In [233]: print regex.match(text)
None
```

Relatedly, `sub` will return a new string with occurrences of the pattern replaced by the a new string:

```
In [234]: print regex.sub('REDACTED', text)
Dave REDACTED
Steve REDACTED
Rob REDACTED
Ryan REDACTED
```

Suppose you wanted to find email addresses and simultaneously segment each address into its 3 components: username, domain name, and domain suffix. To do this, put parentheses around the parts of the pattern to segment:

```
In [235]: pattern = r'([A-Z0-9._%+-]+)@([A-Z0-9.-]+)\.([A-Z]{2,4})'

In [236]: regex = re.compile(pattern, flags=re.IGNORECASE)
```

A match object produced by this modified regex returns a tuple of the pattern components with its `groups` method:

```
In [237]: m = regex.match('wesm@bright.net')

In [238]: m.groups()
Out[238]: ('wesm', 'bright', 'net')
```

`findall` returns a list of tuples when the pattern has groups:

```
In [239]: regex.findall(text)
Out[239]:
[('dave', 'google', 'com'),
 ('steve', 'gmail', 'com'),
 ('rob', 'gmail', 'com'),
 ('ryan', 'yahoo', 'com')]
```

`sub` also has access to groups in each match using special symbols like `\1`, `\2`, etc.:

```
In [240]: print regex.sub(r'Username: \1, Domain: \2, Suffix: \3', text)
Dave Username: dave, Domain: google, Suffix: com
Steve Username: steve, Domain: gmail, Suffix: com
Rob Username: rob, Domain: gmail, Suffix: com
Ryan Username: ryan, Domain: yahoo, Suffix: com
```

There is much more to regular expressions in Python, most of which is outside the book's scope. To give you a flavor, one variation on the above email regex gives names to the match groups:

```
regex = re.compile(r"""
    (?P<username>[A-Z0-9._%+-]+)
    @
    (?P<domain>[A-Z0-9.-]+)
    \.
    (?P<suffix>[A-Z]{2,4})""", flags=re.IGNORECASE|re.VERBOSE)
```

The match object produced by such a regex can produce a handy dict with the specified group names:

```
In [242]: m = regex.match('wesm@bright.net')

In [243]: m.groupdict()
Out[243]: {'domain': 'bright', 'suffix': 'net', 'username': 'wesm'}
```

Table 7-4. Regular expression methods

Argument	Description
findall, finditer	Return all non-overlapping matching patterns in a string. <code>.findall</code> returns a list of all patterns while <code>finditer</code> returns them one by one from an iterator.
match	Match pattern at start of string and optionally segment pattern components into groups. If the pattern matches, returns a match object, otherwise None.
search	Scan string for match to pattern; returning a match object if so. Unlike <code>match</code> , the match can be anywhere in the string as opposed to only at the beginning.
split	Break string into pieces at each occurrence of pattern.
sub, subn	Replace all ( <code>sub</code> ) or first <code>n</code> occurrences ( <code>subn</code> ) of pattern in string with replacement expression. Use symbols <code>\1</code> , <code>\2</code> , ... to refer to match group elements in the replacement string.

## Vectorized string functions in pandas

Cleaning up a messy data set for analysis often requires a lot of string munging and regularization. To complicate matters, a column containing strings will sometimes have missing data:

```
In [244]: data = {'Dave': 'dave@google.com', 'Steve': 'steve@gmail.com',
.....:           'Rob': 'rob@gmail.com', 'Wes': np.nan}

In [245]: data = Series(data)

In [246]: data
Out[246]:
Dave      dave@google.com
Rob       rob@gmail.com
Steve     steve@gmail.com
Wes        NaN
```

```
In [247]: data.isnull()
Out[247]:
Dave      False
Rob      False
Steve     False
Wes       True
```

String and regular expression methods can be applied (passing a `lambda` or other function) to each value using `data.map`, but it will fail on the NA. To cope with this, Series has concise methods for string operations that skip NA values. These are accessed through Series's `str` attribute; for example, we could check whether each email address has 'gmail' in it with `str.contains`:

```
In [248]: data.str.contains('gmail')
Out[248]:
Dave      False
Rob       True
Steve     True
Wes      NaN
```

Regular expressions can be used, too, along with any `re` options like `IGNORECASE`:

```
In [249]: pattern
Out[249]: '([A-Zo-9._%+-]+)@([A-Zo-9.-]+)\\.( [A-Z]{2,4})'

In [250]: data.str.findall(pattern, flags=re.IGNORECASE)
Out[250]:
Dave      [('dave', 'google', 'com')]
Rob      [('rob', 'gmail', 'com')]
Steve    [('steve', 'gmail', 'com')]
Wes      NaN
```

There are a couple of ways to do vectorized element retrieval. Either use `str.get` or index into the `str` attribute:

```
In [251]: matches = data.str.match(pattern, flags=re.IGNORECASE)
```

```
In [252]: matches
Out[252]:
Dave      ('dave', 'google', 'com')
Rob      ('rob', 'gmail', 'com')
Steve    ('steve', 'gmail', 'com')
Wes      NaN
```

```
In [253]: matches.str.get(1)      In [254]: matches.str[0]
Out[253]:                               Out[254]:
Dave      google                  Dave      dave
Rob      gmail                   Rob      rob
Steve    gmail                   Steve    steve
Wes      NaN                     Wes      NaN
```

You can similarly slice strings using this syntax:

```
In [255]: data.str[:5]
Out[255]:
Dave      dave@
Rob      rob@g
Steve    steve
Wes      NaN
```

Table 7-5. Vectorized string methods

Method	Description
cat	Concatenate strings element-wise with optional delimiter
contains	Return boolean array if each string contains pattern/regex
count	Count occurrences of pattern
endswith, startswith	Equivalent to <code>x.endswith(pattern)</code> or <code>x.startswith(pattern)</code> for each element.
findall	Compute list of all occurrences of pattern/regex for each string
get	Index into each element (retrieve i-th element)
join	Join strings in each element of the Series with passed separator
len	Compute length of each string
lower, upper	Convert cases; equivalent to <code>x.lower()</code> or <code>x.upper()</code> for each element.
match	Use <code>re.match</code> with the passed regular expression on each element, returning matched groups as list.
pad	Add whitespace to left, right, or both sides of strings
center	Equivalent to <code>pad(side='both')</code>
repeat	Duplicate values; for example <code>s.str.repeat(3)</code> equivalent to <code>x * 3</code> for each string.
replace	Replace occurrences of pattern/regex with some other string
slice	Slice each string in the Series.
split	Split strings on delimiter or regular expression
strip, rstrip, lstrip	Trim whitespace, including newlines; equivalent to <code>x.strip()</code> (and <code>rstrip</code> , <code>lstrip</code> , respectively) for each element.

## Example: USDA Food Database

The US Department of Agriculture makes available a database of food nutrient information. Ashley Williams, an English hacker, has made available a version of this database in JSON format (<http://ashleyw.co.uk/project/food-nutrient-database>). The records look like this:

```
{
  "id": 21441,
  "description": "KENTUCKY FRIED CHICKEN, Fried Chicken, EXTRA CRISPY,
  Wing, meat and skin with breading",
  "tags": ["KFC"],
  "manufacturer": "Kentucky Fried Chicken",
  "group": "Fast Foods",
  "portions": [
    {
      "amount": 1,
      "unit": "wing, with skin",
      "grams": 68.0
    },
  ]}
```

```

        ],
        ...
    },
    "nutrients": [
        {
            "value": 20.8,
            "units": "g",
            "description": "Protein",
            "group": "Composition"
        },
        ...
    ]
}

```

Each food has a number of identifying attributes along with two lists of nutrients and portion sizes. Having the data in this form is not particularly amenable for analysis, so we need to do some work to wrangle the data into a better form.

After downloading and extracting the data from the link above, you can load it into Python with any JSON library of your choosing. I'll use the built-in Python `json` module:

```

In [256]: import json

In [257]: db = json.load(open('ch07/foods-2011-10-03.json'))

In [258]: len(db)
Out[258]: 6636

```

Each entry in `db` is a dict containing all the data for a single food. The '`nutrients`' field is a list of dicts, one for each nutrient:

	In [259]: db[0].keys()	In [260]: db[0]['nutrients'][0]
Out[259]:	[u'portions', u'description', u'tags', u'nutrients', u'group', u'id', u'manufacturer']	{u'description': u'Protein', u'group': u'Composition', u'units': u'g', u'value': 25.18}
In [261]: nutrients = DataFrame(db[0]['nutrients'])		
In [262]: nutrients[:7]		
Out[262]:		
		description      group    units      value
0		Protein      Composition    g      25.18
1	Total lipid (fat)	Composition    g      29.20
2	Carbohydrate, by difference	Composition    g      3.06
3		Ash            Other     g      3.28
4		Energy        Energy    kcal     376.00
5		Water        Composition    g      39.28
6		Energy        Energy    kJ      1573.00

When converting a list of dicts to a DataFrame, we can specify a list of fields to extract. We'll take the food names, group, id, and manufacturer:

```
In [263]: info_keys = ['description', 'group', 'id', 'manufacturer']
```

```
In [264]: info = DataFrame(db, columns=info_keys)
```

```
In [265]: info[:5]
```

```
Out[265]:
```

	description	group	id	manufacturer
0	Cheese, caraway	Dairy and Egg Products	1008	
1	Cheese, cheddar	Dairy and Egg Products	1009	
2	Cheese, edam	Dairy and Egg Products	1018	
3	Cheese, feta	Dairy and Egg Products	1019	
4	Cheese, mozzarella, part skim milk	Dairy and Egg Products	1028	

```
In [266]: info
```

```
Out[266]:
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 6636 entries, 0 to 6635
Data columns:
description    6636 non-null values
group          6636 non-null values
id             6636 non-null values
manufacturer   5195 non-null values
dtypes: int64(1), object(3)
```

You can see the distribution of food groups with `value_counts`:

```
In [267]: pd.value_counts(info.group)[:10]
```

```
Out[267]:
```

Vegetables and Vegetable Products	812
Beef Products	618
Baked Products	496
Breakfast Cereals	403
Legumes and Legume Products	365
Fast Foods	365
Lamb, Veal, and Game Products	345
Sweets	341
Pork Products	328
Fruits and Fruit Juices	328

Now, to do some analysis on all of the nutrient data, it's easiest to assemble the nutrients for each food into a single large table. To do so, we need to take several steps. First, I'll convert each list of food nutrients to a DataFrame, add a column for the food id, and append the DataFrame to a list. Then, these can be concatenated together with `concat`:

```
nutrients = []

for rec in db:
    fnuts = DataFrame(rec['nutrients'])
    fnuts['id'] = rec['id']
    nutrients.append(fnuts)

nutrients = pd.concat(nutrients, ignore_index=True)
```

If all goes well, `nutrients` should look like this:

```
In [269]: nutrients
Out[269]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 389355 entries, 0 to 389354
Data columns:
description    389355 non-null values
group         389355 non-null values
units          389355 non-null values
value          389355 non-null values
id             389355 non-null values
dtypes: float64(1), int64(1), object(3)
```

I noticed that, for whatever reason, there are duplicates in this DataFrame, so it makes things easier to drop them:

```
In [270]: nutrients.duplicated().sum()
Out[270]: 14179

In [271]: nutrients = nutrients.drop_duplicates()
```

Since 'group' and 'description' is in both DataFrame objects, we can rename them to make it clear what is what:

```
In [272]: col_mapping = {'description' : 'food',
.....:           'group'       : 'fgroup'}

In [273]: info = info.rename(columns=col_mapping, copy=False)

In [274]: info
Out[274]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 6636 entries, 0 to 6635
Data columns:
food        6636 non-null values
fgroup      6636 non-null values
id          6636 non-null values
manufacturer 5195 non-null values
dtypes: int64(1), object(3)

In [275]: col_mapping = {'description' : 'nutrient',
.....:           'group'       : 'nutgroup'}
```

```
In [276]: nutrients = nutrients.rename(columns=col_mapping, copy=False)

In [277]: nutrients
Out[277]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 375176 entries, 0 to 389354
Data columns:
nutrient     375176 non-null values
nutgroup     375176 non-null values
units        375176 non-null values
value        375176 non-null values
```

```
id           375176 non-null values  
dtypes: float64(1), int64(1), object(3)
```

With all of this done, we're ready to merge `info` with `nutrients`:

```
In [278]: ndata = pd.merge(nutrients, info, on='id', how='outer')
```

```
In [279]: ndata
```

```
Out[279]:  
<class 'pandas.core.frame.DataFrame'>  
Int64Index: 375176 entries, 0 to 375175  
Data columns:  
nutrient      375176 non-null values  
nutgroup      375176 non-null values  
units         375176 non-null values  
value          375176 non-null values  
id            375176 non-null values  
food          375176 non-null values  
fgroup        375176 non-null values  
manufacturer  293054 non-null values  
dtypes: float64(1), int64(1), object(6)
```

```
In [280]: ndata.ix[30000]
```

```
Out[280]:  
nutrient                  Folic acid  
nutgroup                  Vitamins  
units                     mcg  
value                      0  
id                        5658  
food          Ostrich, top loin, cooked  
fgroup        Poultry Products  
manufacturer  
Name: 30000
```

The tools that you need to slice and dice, aggregate, and visualize this dataset will be explored in detail in the next two chapters, so after you get a handle on those methods you might return to this dataset. For example, we could a plot of median values by food group and nutrient type (see [Figure 7-1](#)):

```
In [281]: result = ndata.groupby(['nutrient', 'fgroup'])['value'].quantile(0.5)
```

```
In [282]: result['Zinc, Zn'].order().plot(kind='barh')
```

With a little cleverness, you can find which food is most dense in each nutrient:

```
by_nutrient = ndata.groupby(['nutgroup', 'nutrient'])

get_maximum = lambda x: x.xs(x.value.idxmax())
get_minimum = lambda x: x.xs(x.value.idxmin())

max_foods = by_nutrient.apply(get_maximum)[['value', 'food']]

# make the food a little smaller
max_foods.food = max_foods.food.str[:50]
```

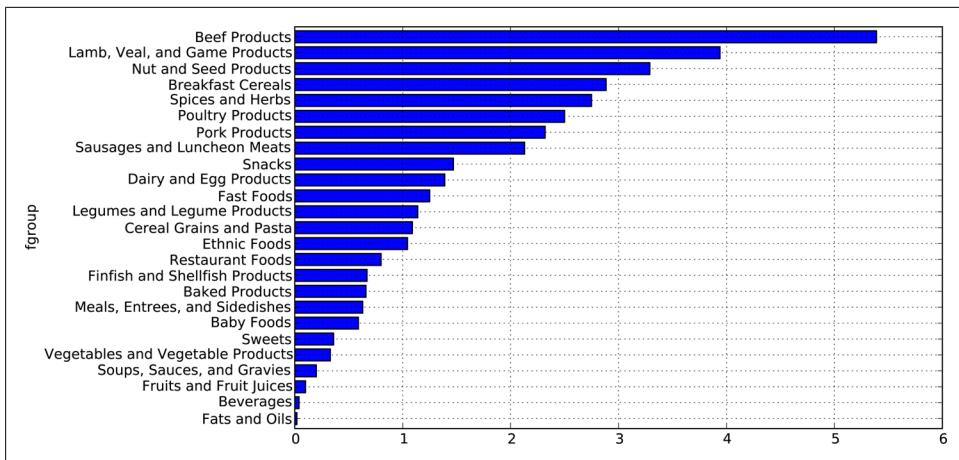


Figure 7-1. Median Zinc values by nutrient group

The resulting DataFrame is a bit too large to display in the book; here is just the 'Amino Acids' nutrient group:

```
In [284]: max_foods.ix['Amino Acids']['food']
Out[284]:
nutrient
Alanine           Gelatins, dry powder, unsweetened
Arginine          Seeds, sesame flour, low-fat
Aspartic acid     Soy protein isolate
Cystine           Seeds, cottonseed flour, low fat (glandless)
Glutamic acid    Soy protein isolate
Glycine           Gelatins, dry powder, unsweetened
Histidine         Whale, beluga, meat, dried (Alaska Native)
Hydroxyproline   KENTUCKY FRIED CHICKEN, Fried Chicken, ORIGINAL R
Isoleucine        Soy protein isolate, PROTEIN TECHNOLOGIES INTERNA
Leucine           Soy protein isolate, PROTEIN TECHNOLOGIES INTERNA
Lysine            Seal, bearded (Oogruk), meat, dried (Alaska Nativ
Methionine       Fish, cod, Atlantic, dried and salted
Phenylalanine    Soy protein isolate, PROTEIN TECHNOLOGIES INTERNA
Proline           Gelatins, dry powder, unsweetened
Serine            Soy protein isolate, PROTEIN TECHNOLOGIES INTERNA
Threonine         Soy protein isolate, PROTEIN TECHNOLOGIES INTERNA
Tryptophan        Sea lion, Steller, meat with fat (Alaska Native)
Tyrosine          Soy protein isolate, PROTEIN TECHNOLOGIES INTERNA
Valine            Soy protein isolate, PROTEIN TECHNOLOGIES INTERNA
Name: food
```

---

# Index

## Symbols

- ! character, 60, 61, 64
- != operator, 91
- !cmd command, 60
- "two-language" problem, 2–3
- # (hash mark), 388
- \$PATH variable, 10
- % character, 398
  - %a datetime format, 293
  - %A datetime format, 293
  - %alias magic function, 61
  - %automagic magic function, 55
  - %b datetime format, 293
  - %B datetime format, 293
  - %bookmark magic function, 60, 62
  - %c datetime format, 293
  - %cd magic function, 60
  - %cpaste magic function, 51–52, 55
  - %d datetime format, 292
  - %D datetime format, 293
  - %d format character, 398
  - %debug magic function, 54–55, 62
  - %dhist magic function, 60
  - %dirs magic function, 60
  - %env magic function, 60
  - %F datetime format, 293
  - %gui magic function, 57
  - %H datetime format, 292
  - %hist magic function, 55, 59
  - %I datetime format, 292
  - %logstart magic function, 60
  - %logstop magic function, 60
  - %lprun magic function, 70, 72
  - %m datetime format, 292
  - %M datetime format, 292
  - %magic magic function, 55
  - %p datetime format, 293
  - %page magic function, 55
  - %paste magic function, 51, 55
  - %pdb magic function, 54, 63
  - %popd magic function, 60
  - %prun magic function, 55, 70
  - %pushd magic function, 60
  - %pwd magic function, 60
  - %quickref magic function, 55
  - %reset magic function, 55, 59
  - %run magic function, 49–50, 55, 386
  - %S datetime format, 292
  - %s format character, 398
  - %time magic function, 55, 67
  - %timeit magic function, 54, 67, 68
  - %U datetime format, 293
  - %w datetime format, 292
  - %W datetime format, 293
  - %who magic function, 55
  - %whos magic function, 55
  - %who\_ls magic function, 55
  - %x datetime format, 293
  - %X datetime format, 293
  - %xdel magic function, 55, 59
  - %xmode magic function, 54
  - %Y datetime format, 292
  - %y datetime format, 292
  - %z datetime format, 293
- & operator, 91
- \* operator, 105
- + operator, 406, 409
- 2012 Federal Election Commission database example, 278–287

We'd like to hear your suggestions for improving our indexes. Send email to [index@oreilly.com](mailto:index@oreilly.com).

bucketing donation amounts, 283–285  
donation statistics by occupation and employer, 280–283  
donation statistics by state, 285–287  
== operator, 393  
>>> prompt, 386  
? (question mark), 49  
[] (brackets), 406, 408  
\ (backslash), 397  
\_ (underscore), 48, 58  
\_\_ (two underscores), 58  
{ } (braces), 413  
| operator, 91

## A

a file mode, 431  
abs function, 96  
accumulate method, 368  
add method, 95, 130, 417  
add\_patch method, 229  
add\_subplot method, 221  
aggfunc option, 277  
aggregate method, 260, 262  
aggregations, 100  
algorithms for sorting, 375–376  
alignment of data, 330–331  
all method, 101, 368  
alpha argument, 233  
and keyword, 398, 401  
annotating in matplotlib, 228–230  
anonymous functions, 424  
any method, 101, 110, 201

arrays  
boolean arrays, 101  
boolean indexing for, 89–92  
conditional logic as operation, 98–100  
creating, 81–82  
creating PeriodIndex from, 312  
data types for, 83–85  
fancy indexing, 92–93  
file input and output with, 103–105  
saving and loading text files, 104–105  
storing on disk in binary format, 103–104  
finding elements in sorted array, 376–377  
in NumPy, 355–362  
concatenating, 357–359  
c\_ object, 359  
layout of in memory, 356–357  
replicating, 360–361  
reshaping, 355–356  
r\_ object, 359  
saving to file, 379–380  
splitting, 357–359  
subsets for, 361–362  
indexes for, 86–89  
operations between, 85–86  
setting values by broadcasting, 367  
slicing, 86–89  
sorting, 101–102  
statistical methods for, 100  
structured arrays, 370–372  
benefits of, 372  
manipulating, 372  
nested data types, 371–372  
swapping axes in, 93–94  
transposing, 93–94  
unique function, 102–103  
where function, 98–100  
arrow function, 229  
as keyword, 393  
asarray function, 82, 379  
asfreq method, 308, 318  
asof method, 334–336  
astype method, 84, 85  
attributes  
in Python, 391  
starting with underscore, 48  
average method, 136  
ax argument, 233  
axes

broadcasting over, 364–367  
concatenating along, 185–188  
labels for, 226–227  
renaming indexes for, 197–198  
swapping in arrays, 93–94  
AxesSubplot object, 221  
axis argument, 188  
axis method, 138

## B

b file mode, 431  
backslash (\), 397  
bar plots, 235–238  
Basemap object, 245, 246  
.bashrc file, 10  
.bash\_profile file, 9  
bbox\_inches option, 231  
benefits  
    of Python, 2–3  
    glue for code, 2  
    solving "two-language" problem with, 2–3  
    of structured arrays, 372  
beta function, 107  
    defined, 342  
between\_time method, 335  
bfill method, 123  
bin edges, 314  
binary data formats, 171–172  
    HDF5, 171–172  
    Microsoft Excel files, 172  
    storing arrays in, 103–104  
binary moving window functions, 324–325  
binary search of lists, 410  
binary universal functions, 96  
binding  
    defined, 390  
    variables, 425  
binomial function, 107  
bisect module, 410  
bookmarking directories in IPython, 62  
Boolean  
    arrays, 101  
    data type, 84, 398  
    indexing for arrays, 89–92  
bottleneck library, 324  
braces ({}), 413  
brackets ([]), 406, 408  
break keyword, 401

broadcasting, 362–367  
    defined, 86, 360, 362  
    over other axes, 364–367  
    setting array values by, 367  
bucketing, 283–285

## C

calendar module, 290  
casting, 84  
cat method, 156, 212  
Categorical object, 199  
ceil function, 96  
center method, 212  
Chaco, 248  
chisquare function, 107  
chunksize argument, 160, 161  
clearing screen shortcut, 53  
clipboard, executing code from, 50–52  
clock function, 67  
close method, 220, 432  
closures, 425–426  
cmd.exe, 7  
collections module, 416  
colons, 387  
cols option, 277  
columns, grouping on, 256–257  
column\_stack function, 359  
combinations function, 430  
combine\_first method, 177, 189  
combining  
    data sources, 336–338  
    data sources, with overlap, 188–189  
    lists, 409  
commands, 65  
    (see also magic commands)  
    debugger, 65  
    history in IPython, 58–60  
        input and output variables, 58–59  
        logging of, 59–60  
        reusing command history, 58  
        searching for, 53  
comment argument, 160  
comments in Python, 388  
compile method, 208  
complex128 data type, 84  
complex256 data type, 84  
complex64 data type, 84  
concat function, 34, 177, 184, 185, 186, 267,  
    357, 359

concatenating  
    along axis, 185–188  
    arrays, 357–359  
conditional logic as array operation, 98–100  
conferences, 12  
configuring matplotlib, 231–232  
conforming, 122  
contains method, 212  
contiguous memory, 381–382  
continue keyword, 401  
continuous return, 348  
convention argument, 314  
converting  
    between string and datetime, 291–293  
    timestamps to periods, 311  
coordinated universal time (UTC), 303  
copy argument, 181  
copy method, 118  
copysign function, 96  
corr method, 140  
correlation, 139–141  
corrwith method, 140  
cos function, 96  
cosh function, 96  
count method, 139, 206, 212, 261, 407  
Counter class, 21  
cov method, 140  
covariance, 139–141  
CPython, 7  
cross-section, 329  
crosstab function, 277–278  
crowdsourcing, 241  
CSV files, 163–165, 242  
Ctrl-A keyboard shortcut, 53  
Ctrl-B keyboard shortcut, 53  
Ctrl-C keyboard shortcut, 53  
Ctrl-E keyboard shortcut, 53  
Ctrl-F keyboard shortcut, 53  
Ctrl-K keyboard shortcut, 53  
Ctrl-L keyboard shortcut, 53  
Ctrl-N keyboard shortcut, 53  
Ctrl-P keyboard shortcut, 53  
Ctrl-R keyboard shortcut, 53  
Ctrl-Shift-V keyboard shortcut, 53  
Ctrl-U keyboard shortcut, 53  
cummax method, 139  
cummin method, 139  
cumprod method, 100, 139  
cumsum method, 100, 139

cumulative returns, 338–340  
currying, 427  
cursor, moving with keyboard, 53  
custom universal functions, 370  
cut function, 199, 200, 201, 268, 283  
Cython project, 2, 382–383  
c\_ object, 359

## D

data aggregation, 259–264  
    returning data in unindexed form, 264  
    using multiple functions, 262–264  
data alignment, 128–132  
    arithmetic methods with fill values, 129–130  
    operations between DataFrame and Series, 130–132  
data munging, 329–340  
    asof method, 334–336  
    combining data, 336–338  
    for data alignment, 330–331  
    for specialized frequencies, 332–334  
data structures for pandas, 112–121  
    DataFrame, 115–120  
    Index objects, 120–121  
    Panel, 152–154  
    Series, 112–115  
data types  
    for arrays, 83–85  
    for ndarray, 83–85  
    for NumPy, 353–354  
        hierarchy of, 354  
    for Python, 395–400  
        boolean data type, 398  
        dates and times, 399–400  
        None data type, 399  
        numeric data types, 395–396  
        str data type, 396–398  
        type casting in, 399  
    for time series data, 290–293  
        converting between string and datetime, 291–293  
    nested, 371–372  
data wrangling  
    manipulating strings, 205–211  
        methods for, 206–207  
        vectorized string methods, 210–211  
        with regular expressions, 207–210  
merging data, 177–189

combining data with overlap, 188–189  
concatenating along axis, 185–188  
DataFrame merges, 178–181  
    on index, 182–184  
pivoting, 192–193  
reshaping, 190–191  
transforming data, 194–205  
    discretization, 199–201  
    dummy variables, 203–205  
    filtering outliers, 201–202  
    mapping, 195–196  
    permutation, 202  
    removing duplicates, 194–195  
    renaming axis indexes, 197–198  
    replacing values, 196–197  
USDA food database example, 212–217

databases  
    reading and writing to, 174–176

DataFrame data structure, 22, 27, 112, 115–120  
    arithmetic operations between Series and, 130–132  
    hierarchical indexing using, 150–151  
    merging data with, 178–181

dates and times, 291  
    (see also time series data)  
    data types for, 291, 399–400  
    date ranges, 298  
    datetime type, 291–293, 395, 399  
    DatetimeIndex Index object, 121  
    dateutil package, 291  
    date\_parser argument, 160  
    date\_range function, 298

dayfirst argument, 160

debug function, 66

debugger, IPython  
    in IPython, 62–66

def keyword, 420

defaults  
    profiles, 77  
    values for dicts, 415–416

del keyword, 59, 118, 414

delete method, 122

delimited formats, 163–165

density plots, 238–239

describe method, 138, 243, 267

design tips, 74–76  
    flat is better than nested, 75  
    keeping relevant objects and data alive, 75

overcoming fear of longer files, 75–76

det function, 106

development tools in IPython, 62–72  
    debugger, 62–66  
    profiling code, 68–70  
    profiling function line-by-line, 70–72  
    timing code, 67–68

diag function, 106

dicts, 413–416  
    creating, 415  
    default values for, 415–416  
    dict comprehensions, 418–420  
    grouping on, 257–258  
    keys for, 416  
    returning system environment variables as, 60

diff method, 122, 139

difference method, 417

digitize function, 377

directories  
    bookmarking in IPython, 62  
    changing, commands for, 60

discretization, 199–201

div method, 130

divide function, 96

.dmg file, 9

donation statistics  
    by occupation and employer, 280–283  
    by state, 285–287

dot function, 105, 106, 377

doublequote option, 165

downsampling, 312

dpi (dots-per-inch) option, 231

reload function, 74

drop method, 122, 125

dropna method, 143

drop\_duplicates method, 194

dsplit function, 359

dstack function, 359

dtype object (see data types)

“duck” typing in Python, 392

dummy variables, 203–205

dumps function, 165

duplicated method, 194

duplicates  
    indices, 296–297  
    removing from data, 194–195

dynamically-generated functions, 425

## E

edgecolor option, 231  
edit-compile-run workflow, 45  
eig function, 106  
elif blocks (see if statements)  
else block (see if statements)  
empty function, 82, 83  
empty namespace, 50  
encoding argument, 160  
endswith method, 207, 212  
enumerate function, 412  
environment variables, 8, 60  
EPD (Enthought Python Distribution), 7–9  
equal function, 96  
escapechar option, 165  
ewma function, 323  
ewmcov function, 323  
ewmcov function, 323  
ewmstd function, 323  
ewmvar function, 323  
ExcelFile class, 172  
except block, 403  
exceptions  
    automatically entering debugger after, 55  
    defined, 402  
    handling in Python, 402–404  
exec keyword, 59  
execute-explore workflow, 45  
execution time  
    of code, 55  
    of single statement, 55  
exit command, 386  
exp function, 96  
expanding window mean, 322  
exponentially-weighted functions, 324  
extend method, 409  
extensible markup language (XML) files, 169–170  
eye function, 83

## F

fabs function, 96  
facecolor option, 231  
factor analysis, 342–343  
Factor object, 269  
factors, 342  
fancy indexing  
    defined, 361

    for arrays, 92–93  
ffill method, 123  
figsize argument, 234  
Figure object, 220, 223  
file input/output  
    binary data formats for, 171–172  
        HDF5, 171–172  
        Microsoft Excel files, 172  
    for arrays, 103–105  
        HDF5, 380  
        memory-mapped files, 379–380  
        saving and loading text files, 104–105  
        storing on disk in binary format, 103–104  
    in Python, 430–431  
    saving plot to file, 231  
text files, 155–170  
    delimited formats, 163–165  
    HTML files, 166–170  
    JSON data, 165–166  
    lxml library, 166–170  
    reading in pieces, 160–162  
    writing to, 162–163  
    XML files, 169–170  
with databases, 174–176  
with Web APIs, 173–174  
filling in missing data, 145–146, 270–271  
fillna method, 22, 143, 145, 146, 196, 270, 317  
fill\_method argument, 313  
fill\_value option, 277  
filtering  
    in pandas, 125–128  
    missing data, 143–144  
    outliers, 201–202  
financial applications  
    cumulative returns, 338–340  
    data munging, 329–340  
        asof method, 334–336  
        combining data, 336–338  
        for data alignment, 330–331  
        for specialized frequencies, 332–334  
future contract rolling, 347–350  
grouping for, 340–345  
    factor analysis with, 342–343  
    quartile analysis, 343–345  
linear regression, 350–351  
return indexes, 338–340  
rolling correlation, 350–351

signal frontier analysis, 345–347  
find method, 206, 207  
findall method, 167, 208, 210, 212  
finditer method, 210  
first crossing time, 109  
first method, 136, 261  
flat is better than nested, 75  
flattening, 356  
float data type, 83, 354, 395, 396, 399  
float function, 402  
float128 data type, 84  
float16 data type, 84  
float32 data type, 84  
float64 data type, 84  
floor function, 96  
floor\_divide function, 96  
flow control, 400–405  
    exception handling, 402–404  
    for loops, 401–402  
    if statements, 400–401  
    pass statements, 402  
    range function, 404–405  
    ternary expressions, 405  
    while loops, 402  
    xrange function, 404–405  
flush method, 432  
fmax function, 96  
fmin function, 96  
fname option, 231  
for loops, 85, 100, 401–402, 418, 419  
format option, 231  
frequencies, 299–301  
    converting, 308  
    specialized frequencies, 332–334  
    week of month dates, 301  
frompyfunc function, 370  
from\_csv method, 163  
functions, 389, 420–430  
    anonymous functions, 424  
    are objects, 422–423  
    closures, 425–426  
    currying of, 427  
    extended call syntax for, 426  
    lambda functions, 424  
    namespaces for, 420–421  
    parsing in pandas, 155  
    returning multiple values from, 422  
    scope of, 420–421  
functools module, 427

future contract rolling, 347–350  
futures, 347

## G

gamma function, 107  
gcc command, 9, 11  
generators, 427–430  
    defined, 428  
    generator expressions, 429  
    itertools module for, 429–430  
get method, 167, 172, 212, 415  
getattr function, 391  
get\_chunk method, 162  
get\_dummies function, 203, 205  
get\_value method, 128  
get\_xlim method, 226  
GIL (global interpreter lock), 3  
global scope, 420, 421  
glue for code  
    Python as, 2  
.gov domain, 17  
Granger, Brian, 72  
graphics  
    Chaco, 248  
    mayavi, 248  
greater function, 96  
greater\_equal function, 96  
grid argument, 234  
group keys, 268  
groupby method, 39, 252–259, 297, 316, 343, 377, 429  
    iterating over groups, 255–256  
    on column, 256–257  
    on dict, 257–258  
    on levels, 259  
    resampling with, 316  
    using functions with, 258–259  
    with Series, 257–258  
grouping  
    2012 Federal Election Commission database  
        example, 278–287  
    bucketing donation amounts, 283–285  
    donation statistics by occupation and  
        employer, 280–283  
    donation statistics by state, 285–287  
    apply method, 266–268  
    data aggregation, 259–264  
        returning data in unindexed form, 264  
        using multiple functions, 262–264

filling missing values with group-specific values, 270–271  
for financial applications, 340–345  
factor analysis with, 342–343  
quartile analysis, 343–345  
group weighted average, 273–274  
groupby method, 252–259  
    iterating over groups, 255–256  
    on column, 256–257  
    on dict, 257–258  
    on levels, 259  
    using functions with, 258–259  
    with Series, 257–258  
linear regression for, 274–275  
pivot tables, 275–278  
    cross-tabulation, 277–278  
quantile analysis with, 268–269  
random sampling with, 271–272

## H

Haiti earthquake crisis data example, 241–247  
half-open, 314  
hasattr function, 391  
hash mark (#), 388  
hashability, 416  
HDF5 (hierarchical data format), 171–172,  
    380  
HDFStore class, 171  
header argument, 160  
heapsort sorting method, 376  
hierarchical data format (HDF5), 171–172,  
    380  
hierarchical indexing  
    in pandas, 147–151  
        sorting levels, 149–150  
        summary statistics by level, 150  
        with DataFrame columns, 150–151  
    reshaping data with, 190–191  
hist method, 238  
histograms, 238–239  
history of commands, searching, 53  
homogeneous data container, 370  
how argument, 181, 313, 316  
hsplit function, 359  
hstack function, 358  
HTML files, 166–170  
HTML Notebook in IPython, 72  
Hunter, John D., 5, 219  
hyperbolic trigonometric functions, 96

## I

icol method, 128, 152  
IDEs (Integrated Development Environments),  
    11, 52  
idxmax method, 138  
idxmin method, 138  
if statements, 400–401, 415  
ifilter function, 430  
iget\_value method, 152  
ignore\_index argument, 188  
imap function, 430  
import directive  
    in Python, 392–393  
    usage of in this book, 13  
imshow function, 98  
in keyword, 409  
in-place sort, 373  
in1d method, 103  
indentation  
    in Python, 387–388  
    IndentationError event, 51  
index method, 206, 207  
Index objects data structure, 120–121  
indexes  
    defined, 112  
    for arrays, 86–89  
    for axis, 197–198  
    for TimeSeries class, 294–296  
    hierarchical indexing, 147–151  
        reshaping data with, 190–191  
        sorting levels, 149–150  
        summary statistics by level, 150  
        with DataFrame columns, 150–151  
    in pandas, 136  
    integer indexing, 151–152  
    merging data on, 182–184  
index\_col argument, 160  
indirect sorts, 374–375, 374  
input variables, 58–59  
insert method, 122, 408  
iinsert method, 410  
int data type, 83, 395, 399  
int16 data type, 84  
int32 data type, 84  
int64 data type, 84  
Int64Index Index object, 121  
int8 data type, 84  
integer arrays, indexing using (see fancy indexing)

integer indexing, 151–152  
Integrated Development Environments (IDEs),  
    11, 52  
interpreted languages  
    defined, 386  
    Python interpreter, 386  
interrupting code, 50, 53  
intersect1d method, 103  
intersection method, 122, 417  
intervals of time, 289  
inv function, 106  
inverse trigonometric functions, 96  
.ipynb files, 72  
IPython, 5  
    bookmarking directories, 62  
    command history in, 58–60  
        input and output variables, 58–59  
        logging of, 59–60  
        reusing command history, 58  
    design tips, 74–76  
        flat is better than nested, 75  
        keeping relevant objects and data alive,  
            75  
        overcoming fear of longer files, 75–76  
    development tools, 62–72  
        debugger, 62–66  
        profiling code, 68–70  
        profiling function line-by-line, 70–72  
        timing code, 67–68  
    executing code from clipboard, 50–52  
    HTML Notebook in, 72  
    integration with IDEs and editors, 52  
    integration with matplotlib, 56–57  
    keyboard shortcuts for, 52  
    magic commands in, 54–55  
    making classes output correctly, 76  
    object introspection in, 48–49  
    profiles for, 77–78  
    Qt console for, 55  
    Quick Reference Card for, 55  
    reloading module dependencies, 74  
    %run command in, 49–50  
    shell commands in, 60–61  
    tab completion in, 47–48  
    tracebacks in, 53–54  
ipython\_config.py file, 77  
irow method, 128, 152  
is keyword, 393  
isdisjoint method, 417

isfinite function, 96  
isin method, 141–142  
isinf function, 96  
isinstance function, 391  
isnull method, 96, 114, 143  
issubdtype function, 354  
issubset method, 417  
issuperset method, 417  
is\_monotonic method, 122  
is\_unique method, 122  
iter function, 392  
iterating over groups, 255–256  
iterator argument, 160  
iterator protocol, 392, 427  
itertools module, 429–430, 429  
ix\_ function, 93

## J

join method, 184, 206, 212  
JSON (JavaScript Object Notation), 18, 165–  
    166, 213

## K

KDE (kernel density estimate) plots, 239  
keep\_date\_col argument, 160  
kernels, 239  
key-value pairs, 413  
keyboard shortcuts, 53  
    for deleting text, 53  
    for IPython, 52  
KeyboardInterrupt event, 50  
keys  
    argument, 188  
    for dicts, 416  
    method, 414  
keyword arguments, 389, 420  
kind argument, 234, 314  
kurt method, 139

## L

label argument, 233, 313, 315  
lambda functions, 211, 262, 424  
last method, 261  
layout of arrays in memory, 356–357  
left argument, 181  
left\_index argument, 181  
left\_on argument, 181  
legends in matplotlib, 228

len function, 212, 258  
less function, 96  
less\_equal function, 96  
level keyword, 259  
levels  
    defined, 147  
    grouping on, 259  
    sorting, 149–150  
    summary statistics by, 150  
lexicographical sort  
    defined, 375  
    lexsort method, 374  
libraries, 3–6  
    IPython, 5  
    matplotlib, 5  
    NumPy, 4  
    pandas, 4–5  
    SciPy, 6  
limit argument, 313  
linalg function, 105  
line plots, 232–235  
linear algebra, 105–106  
linear regression, 274–275, 350–351  
lineterminator option, 164  
line\_profiler extension, 70  
Linux, setting up on, 10–11  
list comprehensions, 418–420  
    nested list comprehensions, 419–420  
list function, 408  
lists, 408–411  
    adding elements to, 408–409  
    binary search of, 410  
    combining, 409  
    insertion into sorted, 410  
    list comprehensions, 418–420  
    removing elements from, 408–409  
    slicing, 410–411  
    sorting, 409–410  
ljust method, 207  
load function, 103, 379  
load method, 171  
loads function, 18  
local scope, 420  
localizing time series data, 304–305  
loffset argument, 313, 316  
log function, 96  
log1p function, 96  
log2 function, 96  
logging command history in IPython, 59–60  
logical\_and function, 96  
logical\_not function, 96  
logical\_or function, 96  
logical\_xor function, 96  
logy argument, 234  
long format, 192  
long type, 395  
longer files overcoming fear of, 75–76  
lower method, 207, 212  
lstrip method, 207, 212  
ltsq function, 106  
lxml library, 166–170

## M

mad method, 139  
magic methods, 48, 54–55  
main function, 75  
manipulating structured arrays, 372  
many-to-many merge, 179  
many-to-one merge, 178  
map method, 133, 195–196, 211, 280, 423  
margins, 275  
markers, 224  
match method, 208–212  
matplotlib, 5, 219–232  
    annotating in, 228–230  
    axis labels in, 226–227  
    configuring, 231–232  
    integrating with IPython, 56–57  
    legends in, 228  
    saving to file, 231  
    styling for, 224–225  
    subplots in, 220–224  
    ticks in, 226–227  
    title in, 226–227  
matplotlibrc file, 232  
matrix operations in NumPy, 377–379  
max method, 101, 136, 139, 261, 428  
maximum function, 95, 96  
mayavi, 248  
mean method, 100, 139, 253, 259, 261, 265  
median method, 139, 261  
memmap object, 379  
memory, layout of arrays in, 356–357  
memory-mapped files  
    defined, 379  
    saving arrays to file, 379–380  
mergesort sorting method, 375, 376  
merging data, 177–189

combining data with overlap, 188–189  
concatenating along axis, 185–188  
DataFrame merges, 178–181  
on index, 182–184  
meshgrid function, 97  
methods  
    defined, 389  
    for tuples, 407  
    in Python, 389  
    starting with underscore, 48  
Microsoft Excel files, 172  
.mil domain, 17  
min method, 101, 136, 139, 261, 428  
minimum function, 96  
missing data, 142–146  
    filling in, 145–146  
    filtering out, 143–144  
mod function, 96  
modf function, 95  
modules, 392  
momentum, 343

**N**

NA data type, 143  
names argument, 160, 188  
namespaces  
    defined, 420  
    in Python, 420–421  
naming trends  
    in US baby names 1880–2010 example, 36–43  
    boy names that became girl names, 42–43  
    measuring increase in diversity, 37–40  
    revolution of last letter, 40–41

NaN (not a number), 101, 114, 143  
na\_values argument, 160  
ncols option, 223  
ndarray, 80  
    Boolean indexing, 89–92  
    creating arrays, 81–82  
    data types for, 83–85  
    fancy indexing, 92–93  
    indexes for, 86–89  
    operations between arrays, 85–86  
    slicing arrays, 86–89  
    swapping axes in, 93–94  
    transposing, 93–94  
nested code, 75  
nested data types, 371–372  
nested list comprehensions, 419–420  
New York MTA (Metropolitan Transportation Authority), 169  
None data type, 395, 399  
normal function, 107, 110  
normalized timestamps, 298  
NoSQL databases, 176  
not a number (NaN), 101, 114, 143  
NotebookCloud, 72  
notnull method, 114, 143  
not\_equal function, 96  
.npy files, 103  
.npz files, 104  
nrows argument, 160, 223  
nuisance column, 254  
numeric data types, 395–396  
NumPy, 4  
    arrays in, 355–362

conditional logic as array operation, 98–100  
methods for boolean arrays, 101  
sorting arrays, 101–102  
statistical methods, 100  
unique function, 102–103  
data types for, 353–354  
file input and output with arrays, 103–105  
    saving and loading text files, 104–105  
    storing on disk in binary format, 103–104  
linear algebra, 105–106  
matrix operations in, 377–379  
ndarray arrays, 80  
    Boolean indexing, 89–92  
    creating, 81–82  
    data types for, 83–85  
    fancy indexing, 92–93  
    indexes for, 86–89  
    operations between arrays, 85–86  
    slicing arrays, 86–89  
    swapping axes in, 93–94  
    transposing, 93–94  
numpy-discussion (mailing list), 12  
performance of, 380–383  
    contiguous memory, 381–382  
    Cython project, 382–383  
random number generation, 106–107  
random walks example, 108–110  
sorting, 373–377  
    algorithms for, 375–376  
    finding elements in sorted array, 376–377  
    indirect sorts, 374–375  
structured arrays in, 370–372  
    benefits of, 372  
    manipulating, 372  
    nested data types, 371–372  
universal functions for, 95–96, 367–370  
    custom, 370  
    in pandas, 132–133  
    instance methods for, 368–369

offsets for time series data, 302–303  
OHLC (Open-High-Low-Close) resampling, 316  
ols function, 351  
Olson database, 303  
on argument, 181  
ones function, 82  
open function, 430  
Open-High-Low-Close (OHLC) resampling, 316  
operators in Python, 393  
or keyword, 401  
order method, 375  
OS X, setting up Python on, 9–10  
outer method, 368, 369  
outliers, filtering, 201–202  
output variables, 58–59

## P

pad method, 212  
pairs plot, 241  
pandas, 4–5  
    arithmetic and data alignment, 128–132  
        arithmetic methods with fill values, 129–130  
    operations between DataFrame and Series, 130–132  
data structures for, 112–121  
    DataFrame, 115–120  
    Index objects, 120–121  
    Panel, 152–154  
    Series, 112–115  
drop function, 125  
filtering in, 125–128  
handling missing data, 142–146  
    filling in, 145–146  
    filtering out, 143–144  
hierarchical indexing in, 147–151  
    sorting levels, 149–150  
    summary statistics by level, 150  
    with DataFrame columns, 150–151  
indexes in, 136  
indexing options, 125–128  
integer indexing, 151–152  
NumPy universal functions with, 132–133  
plotting with, 232  
    bar plots, 235–238  
    density plots, 238–239  
    histograms, 238–239

## 0

object introspection, 48–49  
object model, 388  
object type, 84  
objectify function, 166, 169  
objs argument, 188

line plots, 232–235  
scatter plots, 239–241  
ranking data in, 133–135  
reductions in, 137–142  
reindex function, 122–124  
selecting in objects, 125–128  
sorting in, 133–135  
summary statistics in  
    correlation and covariance, 139–141  
    isin function, 141–142  
    unique function, 141–142  
    value\_counts function, 141–142  
usa.gov data from bit.ly example with, 21–26

Panel data structure, 152–154  
panels, 329  
parse method, 291  
parse\_dates argument, 160  
partial function, 427  
partial indexing, 147  
pass statements, 402  
passing by reference, 390  
pasting  
    keyboard shortcut for, 53  
    magic command for, 55  
patches, 229  
path argument, 160  
Path variable, 8  
pct\_change method, 139  
pdb debugger, 62  
.pdf files, 231  
percentileofscore function, 326  
Pérez, Fernando, 45, 219  
performance  
    and time series data, 327–328  
    of NumPy, 380–383  
    contiguous memory, 381–382  
    Cython project, 382–383  
Period class, 307  
PeriodIndex Index object, 121, 311, 312  
periods, 307–312  
    converting timestamps to, 311  
    creating PeriodIndex from arrays, 312  
    defined, 289, 307  
    frequency conversion for, 308  
    instead of timestamps, 333–334  
    quarterly periods, 309–310  
    resampling with, 318–319  
period\_range function, 307, 310

permutation, 202  
pickle serialization, 170  
pinv function, 106  
pivoting data  
    cross-tabulation, 277–278  
    defined, 189  
    pivot method, 192–193  
    pivot\_table method, 29, 275–278  
pivot\_table aggregation type, 275  
plot method, 23, 36, 41, 220, 224, 232, 239, 246, 319  
plotting  
    Haiti earthquake crisis data example, 241–247  
    time series data, 319–320  
with matplotlib, 219–232  
    annotating in, 228–230  
    axis labels in, 226–227  
    configuring, 231–232  
    legends in, 228  
    saving to file, 231  
    styling for, 224–225  
    subplots in, 220–224  
    ticks in, 226–227  
    title in, 226–227  
with pandas, 232  
    bar plots, 235–238  
    density plots, 238–239  
    histograms, 238–239  
    line plots, 232–235  
    scatter plots, 239–241  
.png files, 231  
pop method, 408, 414  
positional arguments, 389  
power function, 96  
pprint module, 76  
pretty printing  
    and displaying through pager, 55  
    defined, 47  
private attributes, 48  
private methods, 48  
prod method, 261  
profiles  
    defined, 77  
    for IPython, 77–78  
profile\_default directory, 77  
profiling code  
    in IPython, 68–70  
pseudocode, 14

put function, 362  
put method, 362  
.py files, 50, 386, 392  
pydata (Google group), 12  
pylab mode, 219  
pymongo driver, 175  
pyplot module, 220  
pystatsmodels (mailing list), 12  
Python  
    benefits of using, 2–3  
    glue for code, 2  
    solving "two-language" problem with, 2–3  
    data types for, 395–400  
        boolean data type, 398  
        dates and times, 399–400  
        None data type, 399  
        numeric data types, 395–396  
        str data type, 396–398  
        type casting in, 399  
    dict comprehensions in, 418–420  
    dicts in, 413–416  
        creating, 415  
        default values for, 415–416  
        keys for, 416  
    file input/output in, 430–431  
    flow control in, 400–405  
        exception handling, 402–404  
        for loops, 401–402  
        if statements, 400–401  
        pass statements, 402  
        range function, 404–405  
        ternary expressions, 405  
        while loops, 402  
        xrange function, 404–405  
    functions in, 420–430  
        anonymous functions, 424  
        are objects, 422–423  
        closures, 425–426  
        currying of, 427  
        extended call syntax for, 426  
        lambda functions, 424  
        namespaces for, 420–421  
        returning multiple values from, 422  
        scope of, 420–421  
    generators in, 427–430  
        generator expressions, 429  
        itertools module for, 429–430  
IDEs for, 11  
interpreter for, 386  
list comprehensions in, 418–420  
lists in, 408–411  
    adding elements to, 408–409  
    binary search of, 410  
    combining, 409  
    insertion into sorted, 410  
    removing elements from, 408–409  
    slicing, 410–411  
    sorting, 409–410  
Python 2 vs. Python 3, 11  
required libraries, 3–6  
    IPython, 5  
    matplotlib, 5  
    NumPy, 4  
    pandas, 4–5  
    SciPy, 6  
semantics of, 387–395  
    attributes in, 391  
    comments in, 388  
    functions in, 389  
    import directive, 392–393  
    indentation, 387–388  
    methods in, 389  
    mutable objects in, 394–395  
    object model, 388  
    operators for, 393  
    references in, 389–390  
    strict evaluation, 394  
    strongly-typed language, 390–391  
    variables in, 389–390  
    "duck" typing, 392  
sequence functions in, 411–413  
    enumerate function, 412  
    reversed function, 413  
    sorted function, 412  
    zip function, 412–413  
set comprehensions in, 418–420  
sets in, 416–417  
setting up, 6–11  
    on Linux, 10–11  
    on OS X, 9–10  
    on Windows, 7–9  
tuples in, 406–407  
    methods for, 407  
    unpacking, 407  
pytz library, 303

## Q

qcut method, 200, 201, 268, 269, 343  
qr function, 106  
Qt console for IPython, 55  
quantile analysis, 268–269  
quarterly periods, 309–310  
quartile analysis, 343–345  
question mark (?), 49  
quicksort sorting method, 376  
quotechar option, 164  
quoting option, 164

## R

r file mode, 431  
r+ file mode, 431  
Ramachandran, Prabhu, 248  
rand function, 107  
randint function, 107, 202  
randn function, 89, 107  
random number generation, 106–107  
random sampling with grouping, 271–272  
random walks example, 108–110  
range function, 82, 404–405  
ranking data  
    defined, 135  
    in pandas, 133–135  
ravel method, 356, 357  
rc method, 231, 232  
re module, 207  
read method, 432  
read-only mode, 431  
reading  
    from databases, 174–176  
    from text files in pieces, 160–162  
readline functionality, 58  
readlines method, 432  
readshapefile method, 246  
read\_clipboard function, 155  
read\_csv function, 104, 155, 161, 163, 261,  
    430  
read\_frame function, 175  
read\_fwf function, 155  
read\_table function, 104, 155, 158, 163  
recfunctions module, 372  
reduce method, 368, 369  
reduceat method, 369  
reductions, 137  
    (see also aggregations)

defined, 137  
    in pandas, 137–142  
references  
    defined, 389, 390  
    in Python, 389–390  
regress function, 274  
regular expressions (regex)  
    defined, 207  
    manipulating strings with, 207–210  
reindex method, 122–124, 317, 332  
reload function, 74  
remove method, 408, 417  
rename method, 198  
renaming axis indexes, 197–198  
repeat method, 212, 360  
replace method, 196, 206, 212  
replicating arrays, 360–361  
resampling, 312–319, 332  
    defined, 312  
    OHLC (Open-High-Low-Close)  
        resampling, 316  
    upsampling, 316–317  
    with groupby method, 316  
    with periods, 318–319  
reset\_index function, 151  
reshape method, 190–191, 355, 365  
reshaping  
    arrays, 355–356  
    defined, 189  
    with hierarchical indexing, 190–191  
resources, 12  
return statements, 420  
returns  
    cumulative returns, 338–340  
    defined, 338  
    return indexes, 338–340  
reversed function, 413  
rfind method, 207  
right argument, 181  
right\_index argument, 181  
right\_on argument, 181  
rint function, 96  
rjust method, 207  
rollback method, 302  
rollforward method, 302  
rolling, 348  
rolling correlation, 350–351  
rolling\_apply function, 323, 326  
rolling\_corr function, 323, 350

rolling\_count function, 323  
rolling\_cov function, 323  
rolling\_kurt function, 323  
rolling\_mean function, 321, 323  
rolling\_median function, 323  
rolling\_min function, 323  
rolling\_mint function, 323  
rolling\_quantile function, 323, 326  
rolling\_skew function, 323  
rolling\_std function, 323  
rolling\_sum function, 323  
rolling\_var function, 323  
rot argument, 234  
rows option, 277  
row\_stack function, 359  
rstrip method, 207, 212  
r\_object, 359

## S

save function, 103, 379  
save method, 171, 176  
savefig method, 231  
savez function, 104  
saving text files, 104–105  
scatter method, 239  
scatter plots, 239–241  
scatter\_matrix function, 241  
Scientific Python base, 7  
SciPy library, 6  
scipy-user (mailing list), 12  
scope, 420–421  
screen, clearing, 53  
scripting languages, 2  
scripts, 2  
search method, 208, 210  
searchsorted method, 376  
seed function, 107  
seek method, 432  
semantics, 387–395  
    attributes in, 391  
    comments in, 388  
    “duck” typing, 392  
    functions in, 389  
    import directive, 392–393  
    indentation, 387–388  
    methods in, 389  
    mutable objects in, 394–395  
    object model, 388  
    operators for, 393

references in, 389–390  
strict evaluation, 394  
strongly-typed language, 390–391  
variables in, 389–390  
semicolons, 388  
sentinels, 143, 159  
sep argument, 160  
sequence functions, 411–413  
    enumerate function, 412  
    reversed function, 413  
    sorted function, 412  
    zip function, 412–413  
Series data structure, 112–115  
    arithmetic operations between DataFrame  
        and, 130–132  
    grouping with, 257–258  
set comprehensions, 418–420  
set function, 416  
setattr function, 391  
setdefault method, 415  
setdiff1d method, 103  
sets/set comprehensions, 416–417  
setxor1d method, 103  
set\_index function, 151  
set\_index method, 193  
set\_title method, 226  
set\_trace function, 65  
set\_value method, 128  
set\_xlabel method, 226  
set\_xlim method, 226  
set\_xticklabels method, 226  
set\_xticks method, 226  
shapefiles, 246  
shapes, 80, 353  
sharex option, 223, 234  
sharey option, 223, 234  
shell commands in IPython, 60–61  
shifting in time series data, 301–303  
shortcuts, keyboard, 53  
    for deleting text, 53  
    for IPython, 52  
shuffle function, 107  
sign function, 96, 202  
signal frontier analysis, 345–347  
sin function, 96  
sinh function, 96  
size method, 255  
skew method, 139  
skipinitialspace option, 165

skipna method, 138  
skipna option, 137  
skiprows argument, 160  
skip\_footer argument, 160  
slice method, 212  
slicing  
    arrays, 86–89  
    lists, 410–411  
Social Security Administration (SSA), 32  
solve function, 106  
sort argument, 181  
sort method, 101, 373, 409, 424  
sorted function, 412  
sorting  
    arrays, 101–102  
    finding elements in sorted array, 376–377  
    in NumPy, 373–377  
        algorithms for, 375–376  
        finding elements in sorted array, 376–377  
    indirect sorts, 374–375  
    in pandas, 133–135  
        levels, 149–150  
        lists, 409–410  
sortlevel function, 149  
sort\_columns argument, 235  
sort\_index method, 133, 150, 375  
spaces, structuring code with, 387–388  
spacing around subplots, 223–224  
span, 324  
specialized frequencies  
    data munging for, 332–334  
split method, 165, 206, 210, 212, 358  
split-apply-combine, 252  
splitting arrays, 357–359  
SQL databases, 175  
sql module, 175  
SQLite databases, 174  
sqrt function, 95, 96  
square function, 96  
squeeze argument, 160  
SSA (Social Security Administration), 32  
stable sorting, 375  
stacked format, 192  
start index, 411  
startswith method, 207, 212  
statistical methods, 100  
std method, 101, 139, 261  
stdout, 162  
step index, 411  
stop index, 411  
strftime method, 291, 400  
strict evaluation/language, 394  
strides/strided view, 353  
strings  
    converting to datetime, 291–293  
    data types for, 84, 396–398  
    manipulating, 205–211  
        methods for, 206–207  
        vectorized string methods, 210–211  
        with regular expressions, 207–210  
strip method, 207, 212  
strongly-typed languages, 390–391, 390  
strptime method, 291, 400  
structs, 370  
structured arrays, 370–372  
    benefits of, 372  
    defined, 370  
    manipulating, 372  
    nested data types, 371–372  
style argument, 233  
styling for matplotlib, 224–225  
sub method, 130, 209  
subn method, 210  
subperiod, 319  
subplots, 220–224  
subplots method, 222  
subplots\_adjust method, 223  
subplot\_kw option, 223  
subsets for arrays, 361–362  
subtract function, 96  
sudo command, 11  
suffixes argument, 181  
sum method, 100, 132, 137, 139, 259, 261, 330, 428  
summary statistics, 137  
    by level, 150  
    correlation and covariance, 139–141  
    isin function, 141–142  
    unique function, 141–142  
    value\_counts function, 141–142  
superperiod, 319  
svd function, 106  
swapaxes method, 94  
swaplevel function, 149  
swapping axes in arrays, 93–94  
symmetric\_difference method, 417  
syntactic sugar, 14

system commands, defining alias for, 60

## T

tab completion in IPython, 47–48  
tabs, structuring code with, 387–388  
take method, 202, 362  
tan function, 96  
tanh function, 96  
tell method, 432  
terminology, 13–14  
ternary expressions, 405  
text editors, integrating with IPython, 52  
text files, 155–170  
    delimited formats, 163–165  
    HTML files, 166–170  
    JSON data, 165–166  
    lxml library, 166–170  
    reading in pieces, 160–162  
    saving and loading, 104–105  
    writing to, 162–163  
    XML files, 169–170  
TextParser class, 160, 162, 168  
text\_content method, 167  
thousands argument, 160  
thresh argument, 144  
ticks, 226–227  
tile function, 360, 361  
time series data  
    and performance, 327–328  
    data types for, 290–293  
        converting between string and datetime, 291–293  
    date ranges, 298  
    frequencies, 299–301  
        week of month dates, 301  
    moving window functions, 320–326  
        binary moving window functions, 324–325  
        exponentially-weighted functions, 324  
        user-defined, 326  
periods, 307–312  
    converting timestamps to, 311  
    creating PeriodIndex from arrays, 312  
    frequency conversion for, 308  
    quarterly periods, 309–310  
plotting, 319–320  
resampling, 312–319  
    OHLC (Open-High-Low-Close)  
        resampling, 316

upsampling, 316–317  
with groupby method, 316  
with periods, 318–319  
shifting in, 301–303  
    with offsets, 302–303  
time zones in, 303–306  
    localizing objects, 304–305  
    methods for time zone-aware objects, 305–306  
TimeSeries class, 293–297  
    duplicate indices with, 296–297  
    indexes for, 294–296  
    selecting data in, 294–296  
timestamps  
    converting to periods, 311  
    defined, 289  
    using periods instead of, 333–334  
timing code, 67–68  
title in matplotlib, 226–227  
top method, 267, 282  
to\_csv method, 162, 163  
to\_datetime method, 292  
to\_panel method, 154  
to\_period method, 311  
trace function, 106  
tracebacks, 53–54  
transform method, 264–266  
transforming data, 194–205  
    discretization, 199–201  
    dummy variables, 203–205  
    filtering outliers, 201–202  
    mapping, 195–196  
    permutation, 202  
    removing duplicates, 194–195  
    renaming axis indexes, 197–198  
    replacing values, 196–197  
transpose method, 93, 94  
transposing arrays, 93–94  
trellis package, 247  
trigonometric functions, 96  
truncate method, 296  
try/except block, 403, 404  
tuples, 406–407  
    methods for, 407  
    unpacking, 407  
type casting, 399  
type command, 156  
TypeError event, 84, 403  
types, 388

`tz_convert` method, 305  
`tz_localize` method, 304, 305

## U

`U` file mode, 431  
`uint16` data type, 84  
`uint32` data type, 84  
`uint64` data type, 84  
`uint8` data type, 84  
unary functions, 95  
`underscore (_)`, 48, 58  
`unicode` type, 19, 84, 395  
`uniform` function, 107  
`union` method, 103, 122, 204, 417  
`unique` method, 102–103, 122, 141–142, 279  
universal functions, 95–96, 367–370  
    custom, 370  
    in pandas, 132–133  
    instance methods for, 368–369  
universal newline mode, 431  
unpacking tuples, 407  
`unstack` function, 148  
`update` method, 337  
`upper` method, 207, 212  
`upsampling`, 312, 316–317  
US baby names 1880–2010 example, 32–43  
    boy names that became girl names, 42–43  
    measuring increase in diversity, 37–40  
    revolution of last letter, 40–41  
usa.gov data from bit.ly example, 17–26  
USDA (US Department of Agriculture) food  
    database example, 212–217  
`use_index` argument, 234  
UTC (coordinated universal time), 303

## V

`ValueError` event, 402, 403  
`values` method, 414  
`value_counts` method, 141–142  
`var` method, 101, 139, 261  
variables, 55  
    (see also environment variables)  
    deleting, 55  
    displaying, 55  
        in Python, 389–390  
Varoquaux, Gaël, 248  
`vectorization`, 85  
    defined, 97

`vectorize` function, 370  
vectorized string methods, 210–211  
verbose argument, 160  
`verify_integrity` argument, 188  
views, 86, 118  
visualization tools  
    Chaco, 248  
    mayavi, 248  
`vsplit` function, 359  
`vstack` function, 358

## W

`w` file mode, 431  
Wattenberg, Laura, 40  
Web APIs, file input/output with, 173–174  
week of month dates, 301  
when expressions, 394  
where function, 98–100, 188  
while loops, 402  
whitespace, structuring code with, 387–388  
Wickham, Hadley, 252  
Williams, Ashley, 212  
Windows, setting up Python on, 7–9  
working directory  
    changing to passed directory, 60  
    of current system, returning, 60  
wrangling (see data wrangling)  
`write` method, 431  
write-only mode, 431  
`writelines` method, 431  
`writer` method, 165  
writing  
    to databases, 174–176  
    to text files, 162–163

## X

`Xcode`, 9  
`xlim` method, 225, 226  
XML (extensible markup language) files, 169–  
    170  
`xrange` function, 404–405  
`xs` method, 128  
`xticklabels` method, 225

## Y

`yield` keyword, 428  
`ylim` argument, 234  
`ticks` argument, 234