

---

# **pandas: powerful Python data analysis toolkit**

***Release 0.8.1***

**Wes McKinney & PyData Development Team**

July 22, 2012



# CONTENTS

<b>1</b>	<b>What's New</b>	<b>3</b>
1.1	v0.8.1 (July 22, 2012)	3
1.2	v0.8.0 (June 29, 2012)	4
1.3	v0.7.3 (April 12, 2012)	9
1.4	v0.7.2 (March 16, 2012)	13
1.5	v0.7.1 (February 29, 2012)	13
1.6	v0.7.0 (February 9, 2012)	14
1.7	v0.6.1 (December 13, 2011)	19
1.8	v0.6.0 (November 25, 2011)	19
1.9	v0.5.0 (October 24, 2011)	21
1.10	v0.4.3 through v0.4.1 (September 25 - October 9, 2011)	22
<b>2</b>	<b>Installation</b>	<b>25</b>
2.1	Python version support	25
2.2	Binary installers	25
2.3	Dependencies	26
2.4	Optional dependencies	26
2.5	Installing from source	27
2.6	Running the test suite	27
<b>3</b>	<b>Frequently Asked Questions (FAQ)</b>	<b>29</b>
3.1	Migrating from scikits.timeseries to pandas >= 0.8.0	29
<b>4</b>	<b>Package overview</b>	<b>35</b>
4.1	Data structures at a glance	35
4.2	Mutability and copying of data	36
4.3	Getting Support	36
4.4	Credits	36
4.5	Development Team	36
4.6	License	36
<b>5</b>	<b>Intro to Data Structures</b>	<b>39</b>
5.1	Series	39
5.2	DataFrame	43
5.3	Panel	56
<b>6</b>	<b>Essential basic functionality</b>	<b>61</b>
6.1	Head and Tail	61
6.2	Attributes and the raw ndarray(s)	62
6.3	Flexible binary operations	63

6.4	Descriptive statistics . . . . .	67
6.5	Function application . . . . .	72
6.6	Reindexing and altering labels . . . . .	75
6.7	Iteration . . . . .	81
6.8	Vectorized string methods . . . . .	83
6.9	Sorting by index and value . . . . .	86
6.10	Copying, type casting . . . . .	87
6.11	Pickling and serialization . . . . .	88
6.12	Console Output Formatting . . . . .	89
<b>7</b>	<b>Indexing and selecting data</b>	<b>91</b>
7.1	Basics . . . . .	91
7.2	Advanced indexing with labels . . . . .	100
7.3	Index objects . . . . .	104
7.4	Hierarchical indexing (MultiIndex) . . . . .	105
7.5	Adding an index to an existing DataFrame . . . . .	115
7.6	Indexing internal details . . . . .	118
<b>8</b>	<b>Computational tools</b>	<b>119</b>
8.1	Statistical functions . . . . .	119
8.2	Moving (rolling) statistics / moments . . . . .	122
8.3	Exponentially weighted moment functions . . . . .	126
8.4	Linear and panel regression . . . . .	127
<b>9</b>	<b>Working with missing data</b>	<b>135</b>
9.1	Missing data basics . . . . .	135
9.2	Calculations with missing data . . . . .	137
9.3	Cleaning / filling missing data . . . . .	138
9.4	Missing data casting rules and indexing . . . . .	144
<b>10</b>	<b>Group By: split-apply-combine</b>	<b>147</b>
10.1	Splitting an object into groups . . . . .	147
10.2	Iterating through groups . . . . .	151
10.3	Aggregation . . . . .	152
10.4	Transformation . . . . .	155
10.5	Dispatching to instance methods . . . . .	158
10.6	Flexible apply . . . . .	159
10.7	Other useful features . . . . .	160
<b>11</b>	<b>Merge, join, and concatenate</b>	<b>163</b>
11.1	Concatenating objects . . . . .	163
11.2	Database-style DataFrame joining/merging . . . . .	172
<b>12</b>	<b>Reshaping and Pivot Tables</b>	<b>181</b>
12.1	Reshaping by pivoting DataFrame objects . . . . .	181
12.2	Reshaping by stacking and unstacking . . . . .	182
12.3	Reshaping by Melt . . . . .	186
12.4	Combining with stats and GroupBy . . . . .	186
12.5	Pivot tables and cross-tabulations . . . . .	187
12.6	Tiling . . . . .	191
<b>13</b>	<b>Time Series / Date functionality</b>	<b>193</b>
13.1	Time Stamps vs. Time Spans . . . . .	194
13.2	Generating Ranges of Timestamps . . . . .	195
13.3	DateOffset objects . . . . .	198

13.4	Time series-related instance methods . . . . .	203
13.5	Up- and downsampling . . . . .	205
13.6	Time Span Representation . . . . .	207
13.7	Converting between Representations . . . . .	209
13.8	Time Zone Handling . . . . .	211
<b>14</b>	<b>Plotting with matplotlib</b>	<b>215</b>
14.1	Basic plotting: <code>plot</code> . . . . .	215
14.2	Other plotting features . . . . .	222
<b>15</b>	<b>IO Tools (Text, CSV, HDF5, ...)</b>	<b>237</b>
15.1	Clipboard . . . . .	237
15.2	CSV & Text files . . . . .	237
15.3	Excel files . . . . .	250
15.4	HDF5 (PyTables) . . . . .	250
<b>16</b>	<b>Sparse data structures</b>	<b>253</b>
16.1	SparseArray . . . . .	254
16.2	SparseList . . . . .	255
16.3	SparseIndex objects . . . . .	256
<b>17</b>	<b>Caveats and Gotchas</b>	<b>257</b>
17.1	NaN, Integer NA values and NA type promotions . . . . .	257
17.2	Integer indexing . . . . .	259
17.3	Label-based slicing conventions . . . . .	259
17.4	Miscellaneous indexing gotchas . . . . .	260
17.5	Timestamp limitations . . . . .	261
17.6	Parsing Dates from Text Files . . . . .	262
<b>18</b>	<b>rpy2 / R interface</b>	<b>263</b>
18.1	Transferring R data sets into Python . . . . .	263
18.2	Converting DataFrames into R objects . . . . .	264
18.3	Calling R functions with pandas objects . . . . .	264
18.4	High-level interface to R estimators . . . . .	264
<b>19</b>	<b>Related Python libraries</b>	<b>265</b>
19.1	la (larry) . . . . .	265
19.2	scikits.statsmodels . . . . .	265
19.3	scikits.timeseries . . . . .	265
<b>20</b>	<b>Comparison with R / R libraries</b>	<b>267</b>
20.1	<code>data.frame</code> . . . . .	267
20.2	<code>zoo</code> . . . . .	267
20.3	<code>xts</code> . . . . .	267
20.4	<code>plyr</code> . . . . .	267
20.5	<code>reshape / reshape2</code> . . . . .	267
<b>21</b>	<b>API Reference</b>	<b>269</b>
21.1	General functions . . . . .	269
21.2	Series . . . . .	285
21.3	DataFrame . . . . .	308
21.4	Panel . . . . .	344
	<b>Python Module Index</b>	<b>345</b>
	<b>Python Module Index</b>	<b>347</b>



PDF Version **Date:** July 22, 2012 **Version:** 0.8.1

**Binary Installers:** <http://pypi.python.org/pypi/pandas>

**Source Repository:** <http://github.com/pydata/pandas>

**Issues & Ideas:** <https://github.com/pydata/pandas/issues>

**Q&A Support:** <http://stackoverflow.com/questions/tagged/pandas>

**Developer Mailing List:** <http://groups.google.com/group/pystatsmodels>

**pandas** is a **Python** package providing fast, flexible, and expressive data structures designed to make working with “relational” or “labeled” data both easy and intuitive. It aims to be the fundamental high-level building block for doing practical, **real world** data analysis in Python. Additionally, it has the broader goal of becoming **the most powerful and flexible open source data analysis / manipulation tool available in any language**. It is already well on its way toward this goal.

pandas is well suited for many different kinds of data:

- Tabular data with heterogeneously-typed columns, as in an SQL table or Excel spreadsheet
- Ordered and unordered (not necessarily fixed-frequency) time series data.
- Arbitrary matrix data (homogeneously typed or heterogeneous) with row and column labels
- Any other form of observational / statistical data sets. The data actually need not be labeled at all to be placed into a pandas data structure

The two primary data structures of pandas, *Series* (1-dimensional) and *DataFrame* (2-dimensional), handle the vast majority of typical use cases in finance, statistics, social science, and many areas of engineering. For R users, *DataFrame* provides everything that R’s *data.frame* provides and much more. pandas is built on top of **NumPy** and is intended to integrate well within a scientific computing environment with many other 3rd party libraries.

Here are just a few of the things that pandas does well:

- Easy handling of **missing data** (represented as NaN) in floating point as well as non-floating point data
- Size mutability: columns can be **inserted and deleted** from *DataFrame* and higher dimensional objects
- Automatic and explicit **data alignment**: objects can be explicitly aligned to a set of labels, or the user can simply ignore the labels and let *Series*, *DataFrame*, etc. automatically align the data for you in computations
- Powerful, flexible **group by** functionality to perform split-apply-combine operations on data sets, for both aggregating and transforming data
- Make it **easy to convert** ragged, differently-indexed data in other Python and NumPy data structures into *DataFrame* objects
- Intelligent label-based **slicing**, **fancy indexing**, and **subsetting** of large data sets
- Intuitive **merging** and **joining** data sets
- Flexible **reshaping** and pivoting of data sets
- **Hierarchical** labeling of axes (possible to have multiple labels per tick)
- Robust IO tools for loading data from **flat files** (CSV and delimited), Excel files, databases, and saving / loading data from the ultrafast **HDF5 format**
- **Time series**-specific functionality: date range generation and frequency conversion, moving window statistics, moving window linear regressions, date shifting and lagging, etc.

Many of these principles are here to address the shortcomings frequently experienced using other languages / scientific research environments. For data scientists, working with data is typically divided into multiple stages: munging and

cleaning data, analyzing / modeling it, then organizing the results of the analysis into a form suitable for plotting or tabular display. pandas is the ideal tool for all of these tasks.

Some other notes

- pandas is **fast**. Many of the low-level algorithmic bits have been extensively tweaked in [Cython](#) code. However, as with anything else generalization usually sacrifices performance. So if you focus on one feature for your application you may be able to create a faster specialized tool.
- pandas will soon become a dependency of [statsmodels](#), making it an important part of the statistical computing ecosystem in Python.
- pandas has been used extensively in production in financial applications.

---

**Note:** This documentation assumes general familiarity with NumPy. If you haven't used NumPy much or at all, do invest some time in [learning about NumPy](#) first.

---

See the package overview for more detail about what's in the library.



# WHAT'S NEW

These are new features and improvements of note in each release.

## 1.1 v0.8.1 (July 22, 2012)

This release includes a few new features, performance enhancements, and over 30 bug fixes from 0.8.0. New features include notably NA friendly string processing functionality and a series of new plot types and options.

### 1.1.1 New features

- Add *vectorized string processing methods* accessible via `Series.str` (GH620)
- Add option to disable adjustment in EWMA (GH1584)
- *Radviz plot* (GH1566)
- *Parallel coordinates plot*
- *Bootstrap plot*
- Per column styles and secondary y-axis plotting (GH1559)
- New datetime converters millisecond plotting (GH1599)
- Add option to disable “sparse” display of hierarchical indexes (GH1538)
- `Series/DataFrame`’s `set_index` method can *append levels* to an existing `Index/MultiIndex` (GH1569, GH1577)

### 1.1.2 Performance improvements

- Improved implementation of rolling min and max (thanks to Bottleneck !)
- Add accelerated ‘median’ `GroupBy` option (GH1358)
- Significantly improve the performance of parsing ISO8601-format date strings with `DatetimeIndex` or `to_datetime` (GH1571)
- Improve the performance of `GroupBy` on single-key aggregations and use with Categorical types
- Significant datetime parsing performance improvements

## 1.2 v0.8.0 (June 29, 2012)

This is a major release from 0.7.3 and includes extensive work on the time series handling and processing infrastructure as well as a great deal of new functionality throughout the library. It includes over 700 commits from more than 20 distinct authors. Most pandas 0.7.3 and earlier users should not experience any issues upgrading, but due to the migration to the NumPy datetime64 dtype, there may be a number of bugs and incompatibilities lurking. Lingerin incompatibilities will be fixed ASAP in a 0.8.1 release if necessary. See the [full release notes](#) or issue tracker on GitHub for a complete list.

### 1.2.1 Support for non-unique indexes

All objects can now work with non-unique indexes. Data alignment / join operations work according to SQL join semantics (including, if application, index duplication in many-to-many joins)

### 1.2.2 NumPy datetime64 dtype and 1.6 dependency

Time series data are now represented using NumPy’s datetime64 dtype; thus, pandas 0.8.0 now requires at least NumPy 1.6. It has been tested and verified to work with the development version (1.7+) of NumPy as well which includes some significant user-facing API changes. NumPy 1.6 also has a number of bugs having to do with nanosecond resolution data, so I recommend that you steer clear of NumPy 1.6’s datetime64 API functions (though limited as they are) and only interact with this data using the interface that pandas provides.

See the end of the 0.8.0 section for a “porting” guide listing potential issues for users migrating legacy codebases from pandas 0.7 or earlier to 0.8.0.

Bug fixes to the 0.7.x series for legacy NumPy < 1.6 users will be provided as they arise. There will be no more further development in 0.7.x beyond bug fixes.

### 1.2.3 Time series changes and improvements

---

**Note:** With this release, legacy scikits.timeseries users should be able to port their code to use pandas.

---

**Note:** See [documentation](#) for overview of pandas timeseries API.

---

- New datetime64 representation **speeds up join operations and data alignment, reduces memory usage**, and improve serialization / deserialization performance significantly over datetime.datetime
- High performance and flexible **resample** method for converting from high-to-low and low-to-high frequency. Supports interpolation, user-defined aggregation functions, and control over how the intervals and result labeling are defined. A suite of high performance Cython/C-based resampling functions (including Open-High-Low-Close) have also been implemented.
- Revamp of [frequency aliases](#) and support for **frequency shortcuts** like ‘15min’, or ‘1h30min’
- New [DatetimeIndex class](#) supports both fixed frequency and irregular time series. Replaces now deprecated DateRange class
- New PeriodIndex and Period classes for representing [time spans](#) and performing **calendar logic**, including the 12 fiscal quarterly frequencies `<timeseries.quarterly>`. This is a partial port of, and a substantial enhancement to, elements of the scikits.timeseries codebase. Support for conversion between PeriodIndex and DatetimeIndex

- New Timestamp data type subclasses `datetime.datetime`, providing the same interface while enabling working with nanosecond-resolution data. Also provides *easy time zone conversions*.
- Enhanced support for *time zones*. Add `tz_convert` and `tz_localize` methods to `TimeSeries` and `DataFrame`. All timestamps are stored as UTC; Timestamps from `DatetimeIndex` objects with time zone set will be localized to local time. Time zone conversions are therefore essentially free. User needs to know very little about `pytz` library now; only time zone names as strings are required. Time zone-aware timestamps are equal if and only if their UTC timestamps match. Operations between time zone-aware time series with different time zones will result in a UTC-indexed time series.
- Time series **string indexing conveniences** / shortcuts: slice years, year and month, and index values with strings
- Enhanced time series **plotting**; adaptation of `scikits.timeseries` matplotlib-based plotting code
- New `date_range`, `bdate_range`, and `period_range` *factory functions*
- Robust **frequency inference** function `infer_freq` and `inferred_freq` property of `DatetimeIndex`, with option to infer frequency on construction of `DatetimeIndex`
- `to_datetime` function efficiently **parses array of strings** to `DatetimeIndex`. `DatetimeIndex` will parse array or list of strings to `datetime64`
- **Optimized** support for `datetime64`-dtype data in `Series` and `DataFrame` columns
- New NaT (Not-a-Time) type to represent **NA** in timestamp arrays
- Optimize `Series.asof` for looking up “**as of**” values for arrays of timestamps
- Milli, Micro, Nano date offset objects
- Can index time series with `datetime.time` objects to select all data at particular **time of day** (`TimeSeries.at_time`) or **between two times** (`TimeSeries.between_time`)
- Add *tshift* method for leading/lagging using the frequency (if any) of the index, as opposed to a naive lead/lag using `shift`

## 1.2.4 Other new features

- New *cut* and `qcut` functions (like R’s `cut` function) for computing a categorical variable from a continuous variable by binning values either into value-based (`cut`) or quantile-based (`qcut`) bins
- Rename `Factor` to `Categorical` and add a number of usability features
- Add *limit* argument to `fillna/reindex`
- More flexible multiple function application in `GroupBy`, and can pass list (name, function) tuples to get result in particular order with given names
- Add flexible *replace* method for efficiently substituting values
- Enhanced *read\_csv/read\_table* for reading time series data and converting multiple columns to dates
- Add *comments* option to parser functions: `read_csv`, etc.
- Add `:ref`dayfirst`<io.dayfirst>` option to parser functions for parsing international DD/MM/YYYY dates
- Allow the user to specify the CSV reader *dialect* to control quoting etc.
- Handling *thousands* separators in `read_csv` to improve integer parsing.
- Enable unstacking of multiple levels in one shot. Alleviate `pivot_table` bugs (empty columns being introduced)
- Move to `klib`-based hash tables for indexing; better performance and less memory usage than Python’s `dict`

- Add first, last, min, max, and prod optimized GroupBy functions
- New *ordered\_merge* function
- Add flexible *comparison* instance methods eq, ne, lt, gt, etc. to DataFrame, Series
- Improve *scatter\_matrix* plotting function and add histogram or kernel density estimates to diagonal
- Add *'kde'* plot option for density plots
- Support for converting DataFrame to R data.frame through rpy2
- Improved support for complex numbers in Series and DataFrame
- Add *pct\_change* method to all data structures
- Add max\_colwidth configuration option for DataFrame console output
- *Interpolate* Series values using index values
- Can select multiple columns from GroupBy
- Add *update* methods to Series/DataFrame for updating values in place
- Add any and *"all* method to DataFrame

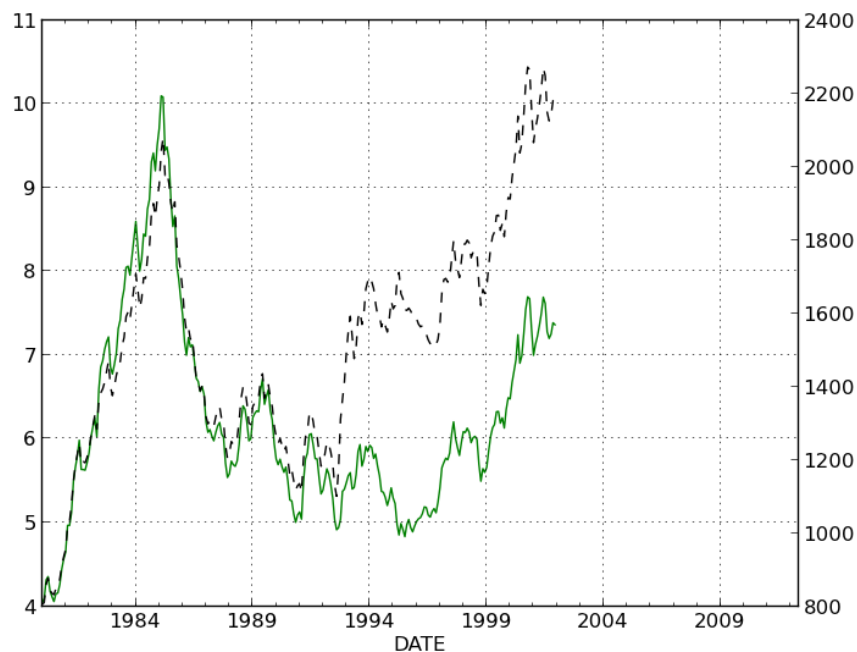
## 1.2.5 New plotting methods

Series.plot now supports a secondary\_y option:

```
In [1]: plt.figure()
Out[1]: <matplotlib.figure.Figure at 0x402f6d0>

In [2]: fx['FR'].plot(style='g')
Out[2]: <matplotlib.axes.AxesSubplot at 0x4052210>

In [3]: fx['IT'].plot(style='k--', secondary_y=True)
Out[3]: <matplotlib.axes.AxesSubplot at 0x4052210>
```



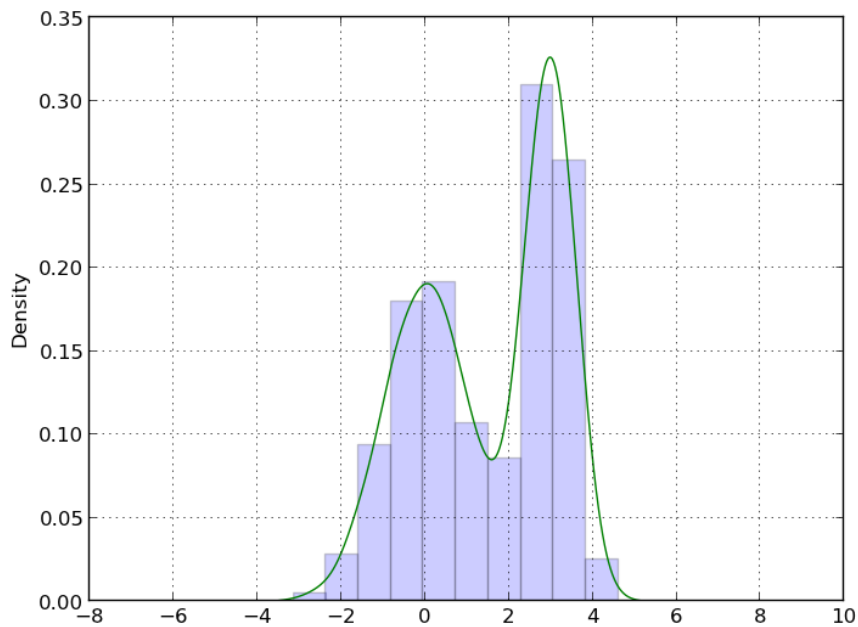
Vytautas Jancauskas, the 2012 GSOC participant, has added many new plot types. For example, 'kde' is a new option:

```
In [4]: s = Series(np.concatenate((np.random.randn(1000),
...:                               np.random.randn(1000) * 0.5 + 3)))
...:
...:
```

```
In [5]: plt.figure()
Out[5]: <matplotlib.figure.Figure at 0x402fcd0>
```

```
In [6]: s.hist(normed=True, alpha=0.2)
Out[6]: <matplotlib.axes.AxesSubplot at 0x4387a50>
```

```
In [7]: s.plot(kind='kde')
Out[7]: <matplotlib.axes.AxesSubplot at 0x4387a50>
```



See [the plotting page](#) for much more.

## 1.2.6 Other API changes

- Deprecation of `offset`, `time_rule`, and `timeRule` arguments names in time series functions. Warnings will be printed until pandas 0.9 or 1.0.

## 1.2.7 Potential porting issues for pandas <= 0.7.3 users

The major change that may affect you in pandas 0.8.0 is that time series indexes use NumPy's `datetime64` data type instead of `dtype=object` arrays of Python's built-in `datetime.datetime` objects. `DateRange` has been replaced by `DatetimeIndex` but otherwise behaved identically. But, if you have code that converts `DateRange` or `Index` objects that used to contain `datetime.datetime` values to plain NumPy arrays, you may have bugs lurking with code using scalar values because you are handing control over to NumPy:

```
In [8]: import datetime
```

```
In [9]: rng = date_range('1/1/2000', periods=10)
```

```
In [10]: rng[5]
Out[10]: <Timestamp: 2000-01-06 00:00:00>

In [11]: isinstance(rng[5], datetime.datetime)
Out[11]: True

In [12]: rng_asarray = np.asarray(rng)

In [13]: scalar_val = rng_asarray[5]

In [14]: type(scalar_val)
Out[14]: numpy.datetime64
```

pandas's `Timestamp` object is a subclass of `datetime.datetime` that has nanosecond support (the nanosecond field store the nanosecond value between 0 and 999). It should substitute directly into any code that used `datetime.datetime` values before. Thus, I recommend not casting `DatetimeIndex` to regular NumPy arrays.

If you have code that requires an array of `datetime.datetime` objects, you have a couple of options. First, the `asobject` property of `DatetimeIndex` produces an array of `Timestamp` objects:

```
In [15]: stamp_array = rng.asobject

In [16]: stamp_array
Out[16]:
Index([2000-01-01 00:00:00, 2000-01-02 00:00:00, 2000-01-03 00:00:00,
       2000-01-04 00:00:00, 2000-01-05 00:00:00, 2000-01-06 00:00:00,
       2000-01-07 00:00:00, 2000-01-08 00:00:00, 2000-01-09 00:00:00,
       2000-01-10 00:00:00], dtype=object)

In [17]: stamp_array[5]
Out[17]: <Timestamp: 2000-01-06 00:00:00>
```

To get an array of proper `datetime.datetime` objects, use the `to_pydatetime` method:

```
In [18]: dt_array = rng.to_pydatetime()

In [19]: dt_array
Out[19]:
array([2000-01-01 00:00:00, 2000-01-02 00:00:00, 2000-01-03 00:00:00,
       2000-01-04 00:00:00, 2000-01-05 00:00:00, 2000-01-06 00:00:00,
       2000-01-07 00:00:00, 2000-01-08 00:00:00, 2000-01-09 00:00:00,
       2000-01-10 00:00:00], dtype=object)

In [20]: dt_array[5]
Out[20]: datetime.datetime(2000, 1, 6, 0, 0)
```

matplotlib knows how to handle `datetime.datetime` but not `Timestamp` objects. While I recommend that you plot time series using `TimeSeries.plot`, you can either use `to_pydatetime` or register a converter for the `Timestamp` type. See [matplotlib documentation](#) for more on this.

**Warning:** There are bugs in the user-facing API with the nanosecond `datetime64` unit in NumPy 1.6. In particular, the string version of the array shows garbage values, and conversion to `dtype=object` is similarly broken.

```
In [21]: rng = date_range('1/1/2000', periods=10)

In [22]: rng
Out[22]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2000-01-01 00:00:00, ..., 2000-01-10 00:00:00]
Length: 10, Freq: D, Timezone: None

In [23]: np.asarray(rng)
Out[23]:
array([1970-01-11 184:00:00, 1970-01-11 208:00:00, 1970-01-11 232:00:00,
       1970-01-11 00:00:00, 1970-01-11 24:00:00, 1970-01-11 48:00:00,
       1970-01-11 72:00:00, 1970-01-11 96:00:00, 1970-01-11 120:00:00,
       1970-01-11 144:00:00], dtype=datetime64[ns])

In [24]: converted = np.asarray(rng, dtype=object)

In [25]: converted[5]
Out[25]: datetime.datetime(1970, 1, 11, 48, 0)
```

**Trust me: don't panic.** If you are using NumPy 1.6 and restrict your interaction with `datetime64` values to pandas's API you will be just fine. There is nothing wrong with the data-type (a 64-bit integer internally); all of the important data processing happens in pandas and is heavily tested. I strongly recommend that you **do not work directly with `datetime64` arrays in NumPy 1.6** and only use the pandas API.

**Support for non-unique indexes:** In the latter case, you may have code inside a `try: ... catch:` block that failed due to the index not being unique. In many cases it will no longer fail (some method like `append` still check for uniqueness unless disabled). However, all is not lost: you can inspect `index.is_unique` and raise an exception explicitly if it is `False` or go to a different code branch.

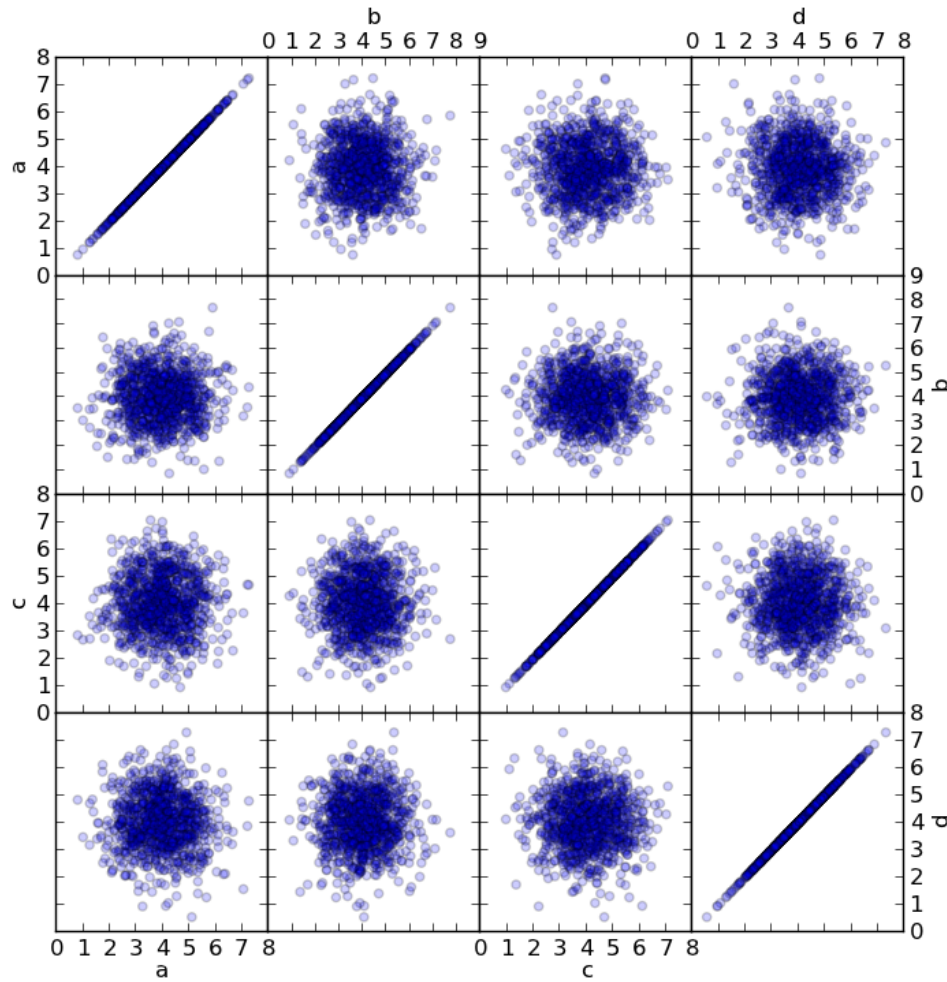
## 1.3 v.0.7.3 (April 12, 2012)

This is a minor release from 0.7.2 and fixes many minor bugs and adds a number of nice new features. There are also a couple of API changes to note; these should not affect very many users, and we are inclined to call them “bug fixes” even though they do constitute a change in behavior. See the [full release notes](#) or issue tracker on GitHub for a complete list.

### 1.3.1 New features

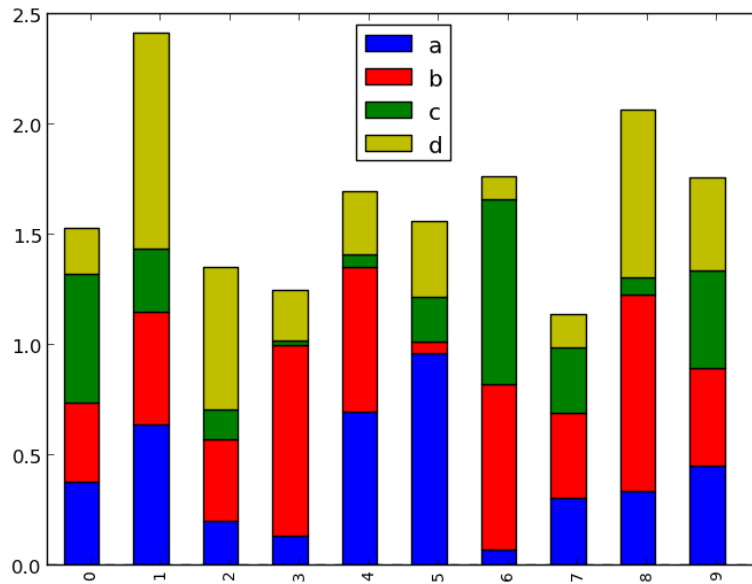
- New *fixed width file reader*, `read_fwf`
- New *scatter\_matrix* function for making a scatter plot matrix

```
from pandas.tools.plotting import scatter_matrix
scatter_matrix(df, alpha=0.2)
```



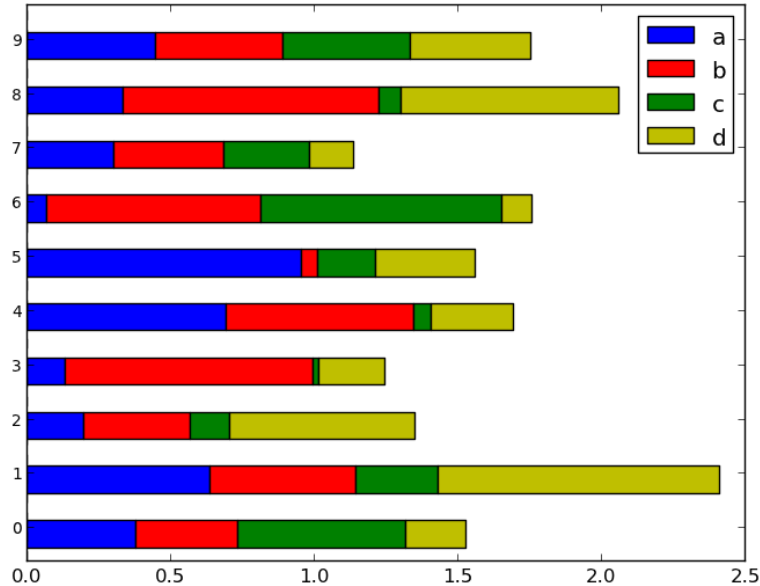
- Add stacked argument to Series and DataFrame's plot method for *stacked bar plots*.

```
df.plot(kind='bar', stacked=True)
```





```
df.plot(kind='barh', stacked=True)
```



- Add log x and y *scaling options* to `DataFrame.plot` and `Series.plot`
- Add kurt methods to `Series` and `DataFrame` for computing kurtosis

### 1.3.2 NA Boolean Comparison API Change

Reverted some changes to how NA values (represented typically as `NaN` or `None`) are handled in non-numeric `Series`:

```
In [26]: series = Series(['Steve', np.nan, 'Joe'])
```

```
In [27]: series == 'Steve'
```

```
Out[27]:
0    True
1   False
2   False
```

```
In [28]: series != 'Steve'
```

```
Out[28]:
0   False
1    True
2    True
```

In comparisons, NA / `NaN` will always come through as `False` except with `!=` which is `True`. *Be very careful* with boolean arithmetic, especially negation, in the presence of NA data. You may wish to add an explicit NA filter into boolean array operations if you are worried about this:

```
In [29]: mask = series == 'Steve'
```

```
In [30]: series[mask & series.notnull()]
```

```
Out[30]: 0    Steve
```

While propagating NA in comparisons may seem like the right behavior to some users (and you could argue on purely technical grounds that this is the right thing to do), the evaluation was made that propagating NA everywhere, including in numerical arrays, would cause a large amount of problems for users. Thus, a “practicality beats purity” approach was taken. This issue may be revisited at some point in the future.

### 1.3.3 Other API Changes

When calling `apply` on a grouped Series, the return value will also be a Series, to be more consistent with the `groupby` behavior with `DataFrame`:

```
In [31]: df = DataFrame({'A' : ['foo', 'bar', 'foo', 'bar',  
.....:                        'foo', 'bar', 'foo', 'foo'],  
.....:                  'B' : ['one', 'one', 'two', 'three',  
.....:                        'two', 'two', 'one', 'three'],  
.....:                  'C' : np.random.randn(8), 'D' : np.random.randn(8)})  
.....:
```

```
In [32]: df
```

```
Out[32]:
```

	A	B	C	D
0	foo	one	1.186498	0.505614
1	bar	one	-0.682439	0.071324
2	foo	two	1.516851	1.546970
3	bar	three	0.347015	-0.562137
4	foo	two	0.099768	-1.572215
5	bar	two	0.601059	0.408282
6	foo	one	-0.539300	0.612042
7	foo	three	0.248445	-0.812354

```
In [33]: grouped = df.groupby('A')['C']
```

```
In [34]: grouped.describe()
```

```
Out[34]:
```

A		
bar	count	3.000000
	mean	0.088545
	std	0.679667
	min	-0.682439
	25%	-0.167712
	50%	0.347015
	75%	0.474037
	max	0.601059
foo	count	5.000000
	mean	0.502452
	std	0.837980
	min	-0.539300
	25%	0.099768
	50%	0.248445
	75%	1.186498
	max	1.516851

```
In [35]: grouped.apply(lambda x: x.order()[-2:]) # top 2 values
```

```
Out[35]:
```

A		
bar	3	0.347015
	5	0.601059
foo	0	1.186498
	2	1.516851

## 1.4 v.0.7.2 (March 16, 2012)

This release targets bugs in 0.7.1, and adds a few minor features.

### 1.4.1 New features

- Add additional tie-breaking methods in `DataFrame.rank` ([GH874](#))
- Add ascending parameter to rank in Series, DataFrame ([GH875](#))
- Add `coerce_float` option to `DataFrame.from_records` ([GH893](#))
- Add `sort_columns` parameter to allow unsorted plots ([GH918](#))
- Enable column access via attributes on `GroupBy` ([GH882](#))
- Can pass dict of values to `DataFrame.fillna` ([GH661](#))
- Can select multiple hierarchical groups by passing list of values in `.ix` ([GH134](#))
- Add `axis` option to `DataFrame.fillna` ([GH174](#))
- Add level keyword to `drop` for dropping values from a level ([GH159](#))

### 1.4.2 Performance improvements

- Use khash for `Series.value_counts`, add `raw` function to `algorithms.py` ([GH861](#))
- Intercept `__builtin__.sum` in `groupby` ([GH885](#))

## 1.5 v.0.7.1 (February 29, 2012)

This release includes a few new features and addresses over a dozen bugs in 0.7.0.

### 1.5.1 New features

- Add `to_clipboard` function to pandas namespace for writing objects to the system clipboard ([GH774](#))
- Add `itertuples` method to `DataFrame` for iterating through the rows of a dataframe as tuples ([GH818](#))
- Add ability to pass `fill_value` and method to `DataFrame` and `Series` `align` method ([GH806](#), [GH807](#))
- Add `fill_value` option to `reindex`, `align` methods ([GH784](#))
- Enable `concat` to produce `DataFrame` from `Series` ([GH787](#))
- Add `between` method to `Series` ([GH802](#))
- Add HTML representation hook to `DataFrame` for the IPython HTML notebook ([GH773](#))
- Support for reading Excel 2007 XML documents using `openpyxl`

### 1.5.2 Performance improvements

- Improve performance and memory usage of `fillna` on `DataFrame`
- Can concatenate a list of `Series` along `axis=1` to obtain a `DataFrame` ([GH787](#))

## 1.6 v.0.7.0 (February 9, 2012)

### 1.6.1 New features

- New unified *merge function* for efficiently performing full gamut of database / relational-algebra operations. Refactored existing join methods to use the new infrastructure, resulting in substantial performance gains (GH220, GH249, GH267)
- New *unified concatenation function* for concatenating Series, DataFrame or Panel objects along an axis. Can form union or intersection of the other axes. Improves performance of `Series.append` and `DataFrame.append` (GH468, GH479, GH273)
- *Can* pass multiple DataFrames to `DataFrame.append` to concatenate (stack) and multiple Series to `Series.append` too
- *Can* pass list of dicts (e.g., a list of JSON objects) to DataFrame constructor (GH526)
- You can now *set multiple columns* in a DataFrame via `__getitem__`, useful for transformation (GH342)
- Handle differently-indexed output values in `DataFrame.apply` (GH498)

```
In [36]: df = DataFrame(randn(10, 4))
```

```
In [37]: df.apply(lambda x: x.describe())
```

```
Out[37]:
```

	0	1	2	3
count	10.000000	10.000000	10.000000	10.000000
mean	0.571888	0.196979	-0.100213	0.058135
std	0.614857	0.676804	0.796945	1.244868
min	-0.684837	-0.584532	-1.881817	-2.867744
25%	0.253154	-0.359561	-0.278064	-0.322515
50%	0.581483	0.084237	0.093037	0.221395
75%	0.960515	0.684092	0.500094	0.763066
max	1.445911	1.477423	0.640984	1.671699

- *Add* `reorder_levels` method to Series and DataFrame (PR534)
- *Add* dict-like `get` function to DataFrame and Panel (PR521)
- *Add* `DataFrame.iterrows` method for efficiently iterating through the rows of a DataFrame
- *Add* `DataFrame.to_panel` with code adapted from `LongPanel.to_long`
- *Add* `reindex_axis` method added to DataFrame
- *Add* `level` option to binary arithmetic functions on DataFrame and Series
- *Add* `level` option to the `reindex` and `align` methods on Series and DataFrame for broadcasting values across a level (GH542, PR552, others)
- *Add* attribute-based item access to Panel and add IPython completion (PR563)
- *Add* `logy` option to `Series.plot` for log-scaling on the Y axis
- *Add* `index` and `header` options to `DataFrame.to_string`
- *Can* pass multiple DataFrames to `DataFrame.join` to join on index (GH115)
- *Can* pass multiple Panels to `Panel.join` (GH115)
- *Added* `justify` argument to `DataFrame.to_string` to allow different alignment of column headers
- *Add* `sort` option to `GroupBy` to allow disabling sorting of the group keys for potential speedups (GH595)

- *Can* pass `MaskedArray` to `Series` constructor ([PR563](#))
- *Add* Panel item access via attributes and IPython completion ([GH554](#))
- Implement `DataFrame.lookup`, fancy-indexing analogue for retrieving values given a sequence of row and column labels ([GH338](#))
- Can pass a *list of functions* to aggregate with `groupby` on a `DataFrame`, yielding an aggregated result with hierarchical columns ([GH166](#))
- Can call `cummin` and `cummax` on `Series` and `DataFrame` to get cumulative minimum and maximum, respectively ([GH647](#))
- `value_range` added as utility function to get min and max of a dataframe ([GH288](#))
- Added encoding argument to `read_csv`, `read_table`, `to_csv` and `from_csv` for non-ascii text ([GH717](#))
- *Added* `abs` method to pandas objects
- *Added* `crosstab` function for easily computing frequency tables
- *Added* `isin` method to index objects
- *Added* `level` argument to `xs` method of `DataFrame`.

## 1.6.2 API Changes to integer indexing

One of the potentially riskiest API changes in 0.7.0, but also one of the most important, was a complete review of how **integer indexes** are handled with regard to label-based indexing. Here is an example:

```
In [38]: s = Series(randn(10), index=range(0, 20, 2))
```

```
In [39]: s
```

```
Out[39]:
0    -0.877467
2     0.814787
4     1.106510
6     0.114919
8    -2.152809
10   -0.518774
12   -0.432823
14   -0.506426
16   -0.471609
18    0.102634
```

```
In [40]: s[0]
```

```
Out[40]: -0.87746706563835053
```

```
In [41]: s[2]
```

```
Out[41]: 0.81478651261064305
```

```
In [42]: s[4]
```

```
Out[42]: 1.1065095963319103
```

This is all exactly identical to the behavior before. However, if you ask for a key **not** contained in the `Series`, in versions 0.6.1 and prior, `Series` would *fall back* on a location-based lookup. This now raises a `KeyError`:

```
In [2]: s[1]
```

```
KeyError: 1
```

This change also has the same impact on DataFrame:

```
In [3]: df = DataFrame(randn(8, 4), index=range(0, 16, 2))
```

```
In [4]: df
      0         1         2         3
0  0.88427  0.3363 -0.1787  0.03162
2  0.14451 -0.1415  0.2504  0.58374
4 -1.44779 -0.9186 -1.4996  0.27163
6 -0.26598 -2.4184 -0.2658  0.11503
8 -0.58776  0.3144 -0.8566  0.61941
10  0.10940 -0.7175 -1.0108  0.47990
12 -1.16919 -0.3087 -0.6049 -0.43544
14 -0.07337  0.3410  0.0424 -0.16037
```

```
In [5]: df.ix[3]
KeyError: 3
```

In order to support purely integer-based indexing, the following methods have been added:

Method	Description
Series.iget_value(i)	Retrieve value stored at location i
Series.iget(i)	Alias for iget_value
DataFrame.irow(i)	Retrieve the i-th row
DataFrame.icol(j)	Retrieve the j-th column
DataFrame.iget_value(i, j)	Retrieve the value at row i and column j

### 1.6.3 API tweaks regarding label-based slicing

Label-based slicing using `ix` now requires that the index be sorted (monotonic) **unless** both the start and endpoint are contained in the index:

```
In [43]: s = Series(randn(6), index=list('gmkaec'))
```

```
In [44]: s
Out[44]:
g  -0.098304
m   0.198159
k  -0.353030
a  -0.076384
e  -0.366746
c  -2.312107
```

Then this is OK:

```
In [45]: s.ix['k':'e']
Out[45]:
k  -0.353030
a  -0.076384
e  -0.366746
```

But this is not:

```
In [12]: s.ix['b':'h']
KeyError 'b'
```

If the index had been sorted, the “range selection” would have been possible:

```
In [46]: s2 = s.sort_index()
```

```
In [47]: s2
```

```
Out[47]:  
a    -0.076384  
c    -2.312107  
e    -0.366746  
g    -0.098304  
k    -0.353030  
m     0.198159
```

```
In [48]: s2.ix['b':'h']
```

```
Out[48]:  
c    -2.312107  
e    -0.366746  
g    -0.098304
```

## 1.6.4 Changes to Series [] operator

As as notational convenience, you can pass a sequence of labels or a label slice to a Series when getting and setting values via [] (i.e. the `__getitem__` and `__setitem__` methods). The behavior will be the same as passing similar input to ix **except in the case of integer indexing**:

```
In [49]: s = Series(randn(6), index=list('acegkm'))
```

```
In [50]: s
```

```
Out[50]:  
a    -0.569813  
c     1.171900  
e    -1.120328  
g     0.388988  
k    -0.290927  
m     0.277157
```

```
In [51]: s[['m', 'a', 'c', 'e']]
```

```
Out[51]:  
m     0.277157  
a    -0.569813  
c     1.171900  
e    -1.120328
```

```
In [52]: s['b':'l']
```

```
Out[52]:  
c     1.171900  
e    -1.120328  
g     0.388988  
k    -0.290927
```

```
In [53]: s['c':'k']
```

```
Out[53]:  
c     1.171900  
e    -1.120328  
g     0.388988  
k    -0.290927
```

In the case of integer indexes, the behavior will be exactly as before (shadowing ndarray):

```
In [54]: s = Series(randn(6), index=range(0, 12, 2))
```

```
In [55]: s[[4, 0, 2]]
```

```
Out[55]:  
4    -0.509155  
0    -0.271227  
2    -1.003010
```

```
In [56]: s[1:5]
```

```
Out[56]:  
2    -1.003010  
4    -0.509155  
6    -0.541238  
8    -0.300009
```

If you wish to do indexing with sequences and slicing on an integer index with label semantics, use `ix`.

## 1.6.5 Other API Changes

- The deprecated `LongPanel` class has been completely removed
- If `Series.sort` is called on a column of a `DataFrame`, an exception will now be raised. Before it was possible to accidentally mutate a `DataFrame`'s column by doing `df[col].sort()` instead of the side-effect free method `df[col].order()` ([GH316](#))
- Miscellaneous renames and deprecations which will (harmlessly) raise `FutureWarning`
- `drop` added as an optional parameter to `DataFrame.reset_index` ([GH699](#))

## 1.6.6 Performance improvements

- *Cythonized GroupBy aggregations* no longer presort the data, thus achieving a significant speedup ([GH93](#)). `GroupBy` aggregations with Python functions significantly sped up by clever manipulation of the `ndarray` data type in Cython ([GH496](#)).
- Better error message in `DataFrame` constructor when passed column labels don't match data ([GH497](#))
- Substantially improve performance of multi-`GroupBy` aggregation when a Python function is passed, reuse `ndarray` object in Cython ([GH496](#))
- Can store objects indexed by tuples and floats in `HDFStore` ([GH492](#))
- Don't print length by default in `Series.to_string`, add *length* option ([GH489](#))
- Improve Cython code for multi-groupby to aggregate without having to sort the data ([GH93](#))
- Improve `MultiIndex` reindexing speed by storing tuples in the `MultiIndex`, test for backwards unpickling compatibility
- Improve column reindexing performance by using specialized Cython take function
- Further performance tweaking of `Series.__getitem__` for standard use cases
- Avoid `Index` dict creation in some cases (i.e. when getting slices, etc.), regression from prior versions
- Friendlier error message in `setup.py` if `NumPy` not installed
- Use common set of NA-handling operations (sum, mean, etc.) in `Panel` class also ([GH536](#))
- Default name assignment when calling `reset_index` on `DataFrame` with a regular (non-hierarchical) index ([GH476](#))



- Use Cythonized groupers when possible in Series/DataFrame stat ops with `level` parameter passed (GH545)
- Ported skiplist data structure to C to speed up `rolling_median` by about 5-10x in most typical use cases (GH374)

## 1.7 v.0.6.1 (December 13, 2011)

### 1.7.1 New features

- Can *append single rows* (as Series) to a DataFrame
- Add Spearman and Kendall rank *correlation* options to `Series.corr` and `DataFrame.corr` (GH428)
- *Added* `get_value` and `set_value` methods to Series, DataFrame, and Panel for very low-overhead access (>2x faster in many cases) to scalar elements (GH437, GH438). `set_value` is capable of producing an enlarged object.
- Add PyQt table widget to sandbox (PR435)
- `DataFrame.align` can *accept Series arguments* and an *axis option* (GH461)
- Implement new *SparseArray* and *SparseList* data structures. `SparseSeries` now derives from `SparseArray` (GH463)
- *Better console printing options* (PR453)
- Implement fast *data ranking* for Series and DataFrame, fast versions of `scipy.stats.rankdata` (GH428)
- Implement *DataFrame.from\_items* alternate constructor (GH444)
- `DataFrame.convert_objects` method for *inferring better dtypes* for object columns (GH302)
- Add *rolling\_corr\_pairwise* function for computing Panel of correlation matrices (GH189)
- Add *margins* option to *pivot\_table* for computing subgroup aggregates (GH114)
- Add `Series.from_csv` function (PR482)
- *Can pass* DataFrame/DataFrame and DataFrame/Series to `rolling_corr/rolling_cov` (GH #462)
- `MultiIndex.get_level_values` can *accept the level name*

### 1.7.2 Performance improvements

- Improve memory usage of `DataFrame.describe` (do not copy data unnecessarily) (PR #425)
- Optimize scalar value lookups in the general case by 25% or more in Series and DataFrame
- Fix performance regression in cross-sectional count in DataFrame, affecting `DataFrame.dropna` speed
- Column deletion in DataFrame copies no data (computes views on blocks) (GH #158)

## 1.8 v.0.6.0 (November 25, 2011)

### 1.8.1 New Features

- *Added* `melt` function to `pandas.core.reshape`
- *Added* `level` parameter to `group by level` in Series and DataFrame descriptive statistics (PR313)

- *Added* `head` and `tail` methods to `Series`, analogous to `DataFrame` (PR296)
- *Added* `Series.isin` function which checks if each value is contained in a passed sequence (GH289)
- *Added* `float_format` option to `Series.to_string`
- *Added* `skip_footer` (GH291) and `converters` (GH343) options to `read_csv` and `read_table`
- *Added* `drop_duplicates` and `duplicated` functions for removing duplicate `DataFrame` rows and checking for duplicate rows, respectively (GH319)
- *Implemented* operators `'&'`, `'|'`, `'^'`, `'-'` on `DataFrame` (GH347)
- *Added* `Series.mad`, mean absolute deviation
- *Added* `QuarterEnd` `DateOffset` (PR321)
- *Added* `dot` to `DataFrame` (GH65)
- *Added* `orient` option to `Panel.from_dict` (GH359, GH301)
- *Added* `orient` option to `DataFrame.from_dict`
- *Added* passing list of tuples or list of lists to `DataFrame.from_records` (GH357)
- *Added* multiple levels to `groupby` (GH103)
- *Allow* multiple columns in `by` argument of `DataFrame.sort_index` (GH92, PR362)
- *Added* `fast_get_value` and `put_value` methods to `DataFrame` (GH360)
- *Added* `cov` instance methods to `Series` and `DataFrame` (GH194, PR362)
- *Added* `kind='bar'` option to `DataFrame.plot` (PR348)
- *Added* `idxmin` and `idxmax` to `Series` and `DataFrame` (PR286)
- *Added* `read_clipboard` function to parse `DataFrame` from clipboard (GH300)
- *Added* `nunique` function to `Series` for counting unique elements (GH297)
- *Made* `DataFrame` constructor use `Series` name if no columns passed (GH373)
- *Support* regular expressions in `read_table/read_csv` (GH364)
- *Added* `DataFrame.to_html` for writing `DataFrame` to HTML (PR387)
- *Added* support for `MaskedArray` data in `DataFrame`, masked values converted to `NaN` (PR396)
- *Added* `DataFrame.boxplot` function (GH368)
- *Can* pass extra args, `kwds` to `DataFrame.apply` (GH376)
- *Implement* `DataFrame.join` with vector on argument (GH312)
- *Added* `legend` boolean flag to `DataFrame.plot` (GH324)
- *Can* pass multiple levels to `stack` and `unstack` (GH370)
- *Can* pass multiple values columns to `pivot_table` (GH381)
- *Use* `Series` name in `GroupBy` for result index (GH363)
- *Added* `raw` option to `DataFrame.apply` for performance if only need `ndarray` (GH309)
- *Added* proper, tested weighted least squares to standard and panel OLS (GH303)

## 1.8.2 Performance Enhancements

- VBENCH Cythonized `cache_readonly`, resulting in substantial micro-performance enhancements throughout the codebase ([GH361](#))
- VBENCH Special Cython matrix iterator for applying arbitrary reduction operations with 3-5x better performance than `np.apply_along_axis` ([GH309](#))
- VBENCH Improved performance of `MultiIndex.from_tuples`
- VBENCH Special Cython matrix iterator for applying arbitrary reduction operations
- VBENCH + DOCUMENT Add `raw` option to `DataFrame.apply` for getting better performance when
- VBENCH Faster cythonized count by level in `Series` and `DataFrame` ([GH341](#))
- VBENCH? Significant `GroupBy` performance enhancement with multiple keys with many “empty” combinations
- VBENCH New Cython vectorized function `map_infer` speeds up `Series.apply` and `Series.map` significantly when passed elementwise Python function, motivated by ([PR355](#))
- VBENCH Significantly improved performance of `Series.order`, which also makes `np.unique` called on a `Series` faster ([GH327](#))
- VBENCH Vastly improved performance of `GroupBy` on axes with a `MultiIndex` ([GH299](#))

## 1.9 v.0.5.0 (October 24, 2011)

### 1.9.1 New Features

- *Added* `DataFrame.align` method with standard join options
- *Added* `parse_dates` option to `read_csv` and `read_table` methods to optionally try to parse dates in the index columns
- *Added* `nrows`, `chunksize`, and `iterator` arguments to `read_csv` and `read_table`. The last two return a new `TextParser` class capable of lazily iterating through chunks of a flat file ([GH242](#))
- *Added* ability to join on multiple columns in `DataFrame.join` ([GH214](#))
- Added private `_get_duplicates` function to `Index` for identifying duplicate values more easily ([ENH5c](#))
- *Added* column attribute access to `DataFrame`.
- *Added* Python tab completion hook for `DataFrame` columns. ([PR233](#), [GH230](#))
- *Implemented* `Series.describe` for `Series` containing objects ([PR241](#))
- *Added* inner join option to `DataFrame.join` when joining on key(s) ([GH248](#))
- *Implemented* selecting `DataFrame` columns by passing a list to `__getitem__` ([GH253](#))
- *Implemented* `&` and `|` to intersect / union `Index` objects, respectively ([GH261](#))
- *Added* `pivot_table` convenience function to pandas namespace ([GH234](#))
- *Implemented* `Panel.rename_axis` function ([GH243](#))
- `DataFrame` will show index level names in console output ([PR334](#))
- *Implemented* `Panel.take`
- *Added* `set_eng_float_format` for alternate `DataFrame` floating point string formatting ([ENH61](#))

- *Added* convenience `set_index` function for creating a `DataFrame` index from its existing columns
- *Implemented* `groupby` hierarchical index level name (GH223)
- *Added* support for different delimiters in `DataFrame.to_csv` (PR244)
- TODO: DOCS ABOUT TAKE METHODS

## 1.9.2 Performance Enhancements

- VBENCH Major performance improvements in file parsing functions `read_csv` and `read_table`
- VBENCH Added Cython function for converting tuples to `ndarray` very fast. Speeds up many `MultiIndex`-related operations
- VBENCH Refactored merging / joining code into a tidy class and disabled unnecessary computations in the float/object case, thus getting about 10% better performance (GH211)
- VBENCH Improved speed of `DataFrame.xs` on mixed-type `DataFrame` objects by about 5x, regression from 0.3.0 (GH215)
- VBENCH With new `DataFrame.align` method, speeding up binary operations between differently-indexed `DataFrame` objects by 10-25%.
- VBENCH Significantly sped up conversion of nested dict into `DataFrame` (GH212)
- VBENCH Significantly speed up `DataFrame.__repr__` and `count` on large mixed-type `DataFrame` objects

## 1.10 v.0.4.3 through v0.4.1 (September 25 - October 9, 2011)

### 1.10.1 New Features

- Added Python 3 support using 2to3 (PR200)
- *Added* name attribute to `Series`, now prints as part of `Series.__repr__`
- *Added* instance methods `isnull` and `notnull` to `Series` (PR209, GH203)
- *Added* `Series.align` method for aligning two series with choice of join method (ENH56)
- *Added* method `get_level_values` to `MultiIndex` (IS188)
- *Set* values in mixed-type `DataFrame` objects via `.ix` indexing attribute (GH135)
- Added new `DataFrame` *methods* `get_dtype_counts` and property `dtypes` (ENHdc)
- Added *ignore\_index* option to `DataFrame.append` to stack `DataFrames` (ENH1b)
- `read_csv` tries to *sniff* delimiters using `csv.Sniffer` (PR146)
- `read_csv` can *read* multiple columns into a `MultiIndex`; `DataFrame`'s `to_csv` method writes out a corresponding `MultiIndex` (PR151)
- `DataFrame.rename` has a new `copy` parameter to *rename* a `DataFrame` in place (ENHed)
- *Enable* unstacking by name (PR142)
- *Enable* `sortlevel` to work by level (PR141)

### 1.10.2 Performance Enhancements

- Altered binary operations on differently-indexed SparseSeries objects to use the integer-based (dense) alignment logic which is faster with a larger number of blocks ([GH205](#))
- Wrote faster Cython data alignment / merging routines resulting in substantial speed increases
- Improved performance of `isnull` and `notnull`, a regression from v0.3.0 ([GH187](#))
- Refactored code related to `DataFrame.join` so that intermediate aligned copies of the data in each `DataFrame` argument do not need to be created. Substantial performance increases result ([GH176](#))
- Substantially improved performance of generic `Index.intersection` and `Index.union`
- Implemented `BlockManager.take` resulting in significantly faster `take` performance on mixed-type `DataFrame` objects ([GH104](#))
- Improved performance of `Series.sort_index`
- Significant groupby performance enhancement: removed unnecessary integrity checks in `DataFrame` internals that were slowing down slicing operations to retrieve groups
- Optimized `_ensure_index` function resulting in performance savings in type-checking `Index` objects
- Wrote fast time series merging / joining methods in Cython. Will be integrated later into `DataFrame.join` and related functions



# INSTALLATION

You have the option to install an [official release](#) or to build the [development version](#). If you choose to install from source and are running Windows, you will have to ensure that you have a compatible C compiler (MinGW or Visual Studio) installed. [How-to install MinGW on Windows](#)

## 2.1 Python version support

Officially Python 2.5 to 2.7 and Python 3.1+, although Python 3 support is less well tested. Python 2.4 support is being phased out since the userbase has shrunk significantly. Continuing Python 2.4 support will require either monetary development support or someone contributing to the project to maintain compatibility.

## 2.2 Binary installers

### 2.2.1 All platforms

Stable installers available on [PyPI](#)

Preliminary builds and installers on the [Pandas download page](#) .

## 2.2.2 Overview

Platform	Distribution	Status	Download / Repository Link	Install method
Windows	all	stable	<i>All platforms</i>	<code>pip install pandas</code>
Mac	all	stable	<i>All platforms</i>	<code>pip install pandas</code>
Linux	Debian	stable	official Debian repository	<code>sudo apt-get install python-pandas</code>
Linux	Debian	unstable (latest packages)	NeuroDebian	<code>sudo apt-get install python-pandas</code>
Linux	Ubuntu	stable	official Ubuntu repository	<code>sudo apt-get install python-pandas</code>
Linux	Ubuntu	unstable (daily builds)	PythonXY PPA; activate by: <code>sudo add-apt-repository ppa:pythonxy/pythonxy-devel</code> && <code>sudo apt-get update</code>	<code>sudo apt-get install python-pandas</code>
Linux	Open-Suse & Fedora	stable	OpenSuse Repository	<code>zypper in python-pandas</code>

## 2.3 Dependencies

- NumPy: 1.6.1 or higher
- python-dateutil 1.5

## 2.4 Optional dependencies

- SciPy: miscellaneous statistical functions
- PyTables: necessary for HDF5-based storage
- matplotlib: for plotting
- `scikits.statsmodels`
  - Needed for parts of `pandas.stats`
- `pytz`
  - Needed for time zone support with `date_range`

---

**Note:** Without the optional dependencies, many useful features will not work. Hence, it is highly recommended that you install these. A packaged distribution like the [Enthought Python Distribution](#) may be worth considering.

---



## 2.5 Installing from source

---

**Note:** Installing from the git repository requires a recent installation of [Cython](#) as the cythonized C sources are no longer checked into source control. Released source distributions will contain the built C files. I recommend installing the latest Cython via `easy_install -U Cython`

---

The source code is hosted at <http://github.com/pydata/pandas>, it can be checked out using git and compiled / installed like so:

```
git clone git://github.com/pydata/pandas.git
cd pandas
python setup.py install
```

On Windows, I suggest installing the MinGW compiler suite following the directions linked to above. Once configured properly, run the following on the command line:

```
python setup.py build --compiler=mingw32
python setup.py install
```

Note that you will not be able to import pandas if you open an interpreter in the source directory unless you build the C extensions in place:

```
python setup.py build_ext --inplace
```

The most recent version of MinGW (any installer dated after 2011-08-03) has removed the ‘-mno-cygwin’ option but Distutils has not yet been updated to reflect that. Thus, you may run into an error like “unrecognized command line option ‘-mno-cygwin’”. Until the bug is fixed in Distutils, you may need to install a slightly older version of MinGW (2011-08-02 installer).

## 2.6 Running the test suite

pandas is equipped with an exhaustive set of unit tests covering about 97% of the codebase as of this writing. To run it on your machine to verify that everything is working (and you have all of the dependencies, soft and hard, installed), make sure you have [nose](#) and run:

```
$ nosetests pandas
.....S.....
.....S.....
.....
-----
Ran 818 tests in 21.631s

OK (SKIP=2)
```



# FREQUENTLY ASKED QUESTIONS (FAQ)

## 3.1 Migrating from `scikits.timeseries` to `pandas` $\geq$ 0.8.0

Starting with `pandas` 0.8.0, users of `scikits.timeseries` should have all of the features that they need to migrate their code to use `pandas`. Portions of the `scikits.timeseries` codebase for implementing calendar logic and timespan frequency conversions (but **not** resampling, that has all been implemented from scratch from the ground up) have been ported to the `pandas` codebase.

The `scikits.timeseries` notions of `Date` and `DateArray` are responsible for implementing calendar logic:

```
In [16]: dt = ts.Date('Q', '1984Q3')

# sic
In [17]: dt
Out[17]: <Q-DEC : 1984Q1>

In [18]: dt.asfreq('D', 'start')
Out[18]: <D : 01-Jan-1984>

In [19]: dt.asfreq('D', 'end')
Out[19]: <D : 31-Mar-1984>

In [20]: dt + 3
Out[20]: <Q-DEC : 1984Q4>
```

`Date` and `DateArray` from `scikits.timeseries` have been reincarnated in `pandas` `Period` and `PeriodIndex`:

```
In [384]: pnow('D') # scikits.timeseries.now()
Out[384]: Period('22-Jul-2012', 'D')

In [385]: Period(year=2007, month=3, day=15, freq='D')
Out[385]: Period('15-Mar-2007', 'D')

In [386]: p = Period('1984Q3')

In [387]: p
Out[387]: Period('1984Q3', 'Q-DEC')

In [388]: p.asfreq('D', 'start')
Out[388]: Period('01-Jul-1984', 'D')
```

```
In [389]: p.asfreq('D', 'end')
Out[389]: Period('30-Sep-1984', 'D')

In [390]: (p + 3).asfreq('T') + 6 * 60 + 30
Out[390]: Period('01-Jul-1985 06:29', 'T')

In [391]: rng = period_range('1990', '2010', freq='A')

In [392]: rng
Out[392]:
<class 'pandas.tseries.period.PeriodIndex'>
freq: A-DEC
[1990, ..., 2010]
length: 21

In [393]: rng.asfreq('B', 'end') - 3
Out[393]:
<class 'pandas.tseries.period.PeriodIndex'>
freq: B
[26-Dec-1990, ..., 28-Dec-2010]
length: 21
```

scikits.timeseries	pandas	Notes
Date	Period	A span of time, from yearly through to secondly
DateArray	PeriodIndex	An array of timespans
convert	resample	Frequency conversion in scikits.timeseries
convert_to_annual	pivot_annual	currently supports up to daily frequency, see <a href="#">:issue:'736'</a>

### 3.1.1 PeriodIndex / DateArray properties and functions

The scikits.timeseries DateArray had a number of information properties. Here are the pandas equivalents:

scikits.timeseries	pandas	Notes
get_steps	np.diff(idx.values)	
has_missing_dates	not idx.is_full	
is_full	idx.is_full	
is_valid	idx.is_monotonic and idx.is_unique	
is_chronological	is_monotonic	
arr.sort_chronologically()	idx.order()	

### 3.1.2 Frequency conversion

Frequency conversion is implemented using the `resample` method on `TimeSeries` and `DataFrame` objects (multiple time series). `resample` also works on panels (3D). Here is some code that resamples daily data to montly with scikits.timeseries:

```
In [394]: import scikits.timeseries as ts

In [395]: data = ts.time_series(np.random.randn(50), start_date='Jan-2000', freq='M')

In [396]: data
Out[396]:
timeseries([ 0.4691 -0.2829 -1.5091 -1.1356  1.2121 -0.1732  0.1192 -1.0442 -0.8618
 -2.1046 -0.4949  1.0718  0.7216 -0.7068 -1.0396  0.2719 -0.425  0.567
  0.2762 -1.0874 -0.6737  0.1136 -1.4784  0.525  0.4047  0.577 -1.715
```

```
-1.0393 -0.3706 -1.1579 -1.3443  0.8449  1.0758 -0.109   1.6436 -1.4694
 0.357  -0.6746 -1.7769 -0.9689 -1.2945  0.4137  0.2767 -0.472  -0.014
-0.3625 -0.0062 -0.9231  0.8957  0.8052],
  dates = [Jan-2012 ... Feb-2016],
  freq   = M)
```

```
In [397]: data.convert('A', func=np.mean)
```

```
Out[397]:
```

```
timeseries([-0.394509620575 -0.24462765889 -0.221632512996 -0.453772693384
 0.8504806638],
  dates = [2012 ... 2016],
  freq   = A-DEC)
```

Here is the equivalent pandas code:

```
In [398]: rng = period_range('Jan-2000', periods=50, freq='M')
```

```
In [399]: data = Series(np.random.randn(50), index=rng)
```

```
In [400]: data
```

```
Out[400]:
```

```
Jan-2000    -1.206412
Feb-2000     2.565646
Mar-2000     1.431256
Apr-2000     1.340309
May-2000    -1.170299
Jun-2000    -0.226169
Jul-2000     0.410835
Aug-2000     0.813850
Sep-2000     0.132003
Oct-2000    -0.827317
Nov-2000    -0.076467
Dec-2000    -1.187678
Jan-2001     1.130127
Feb-2001    -1.436737
Mar-2001    -1.413681
Apr-2001     1.607920
May-2001     1.024180
Jun-2001     0.569605
Jul-2001     0.875906
Aug-2001    -2.211372
Sep-2001     0.974466
Oct-2001    -2.006747
Nov-2001    -0.410001
Dec-2001    -0.078638
Jan-2002     0.545952
Feb-2002    -1.219217
Mar-2002    -1.226825
Apr-2002     0.769804
May-2002    -1.281247
Jun-2002    -0.727707
Jul-2002    -0.121306
Aug-2002    -0.097883
Sep-2002     0.695775
Oct-2002     0.341734
Nov-2002     0.959726
Dec-2002    -1.110336
Jan-2003    -0.619976
```

```
Feb-2003    0.149748
Mar-2003   -0.732339
Apr-2003    0.687738
May-2003    0.176444
Jun-2003    0.403310
Jul-2003   -0.154951
Aug-2003    0.301624
Sep-2003   -2.179861
Oct-2003   -1.369849
Nov-2003   -0.954208
Dec-2003    1.462696
Jan-2004   -1.743161
Feb-2004   -0.826591
Freq: M
```

```
In [401]: data.resample('A', how=np.mean)
```

```
Out [401]:
2000    0.166630
2001   -0.114581
2002   -0.205961
2003   -0.235802
2004   -1.284876
Freq: A-DEC
```

### 3.1.3 Plotting

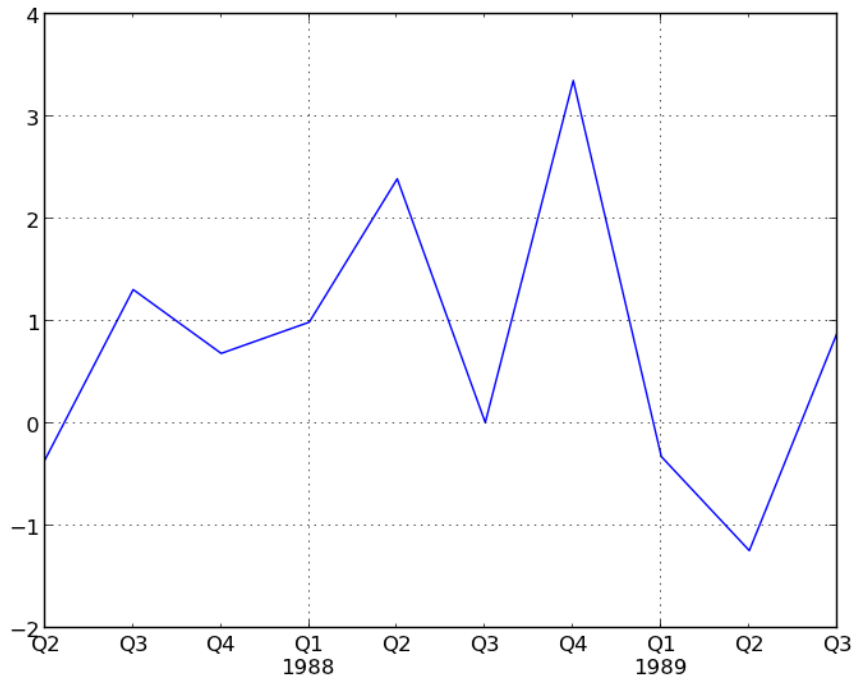
Much of the plotting functionality of `scikits.timeseries` has been ported and adopted to pandas's data structures. For example:

```
In [402]: rng = period_range('1987Q2', periods=10, freq='Q-DEC')
```

```
In [403]: data = Series(np.random.randn(10), index=rng)
```

```
In [404]: plt.figure(); data.plot()
```

```
Out [404]: <matplotlib.axes.AxesSubplot at 0x84e3b50>
```



### 3.1.4 Converting to and from period format

Use the `to_timestamp` and `to_period` instance methods.

### 3.1.5 Treatment of missing data

Unlike `scikits.timeseries`, pandas data structures are not based on NumPy's `MaskedArray` object. Missing data is represented as `NaN` in numerical arrays and either as `None` or `NaN` in non-numerical arrays. Implementing a version of pandas's data structures that use `MaskedArray` is possible but would require the involvement of a dedicated maintainer. Active pandas developers are not interested in this.

### 3.1.6 Resampling with timestamps and periods

`resample` has a `kind` argument which allows you to resample time series with a `DatetimeIndex` to `PeriodIndex`:

```
In [405]: rng = date_range('1/1/2000', periods=200, freq='D')
```

```
In [406]: data = Series(np.random.randn(200), index=rng)
```

```
In [407]: data[:10]
```

```
Out[407]:
2000-01-01    -0.487602
2000-01-02    -0.082240
2000-01-03    -2.182937
2000-01-04     0.380396
2000-01-05     0.084844
2000-01-06     0.432390
2000-01-07     1.519970
2000-01-08    -0.493662
2000-01-09     0.600178
```

```
2000-01-10    0.274230
Freq: D
```

```
In [408]: data.index
Out[408]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2000-01-01 00:00:00, ..., 2000-07-18 00:00:00]
Length: 200, Freq: D, Timezone: None
```

```
In [409]: data.resample('M', kind='period')
Out[409]:
Jan-2000    0.163775
Feb-2000    0.026549
Mar-2000   -0.089563
Apr-2000   -0.079405
May-2000    0.160348
Jun-2000    0.101725
Jul-2000   -0.708770
Freq: M
```

Similarly, resampling from periods to timestamps is possible with an optional interval ('start' or 'end') convention:

```
In [410]: rng = period_range('Jan-2000', periods=50, freq='M')

In [411]: data = Series(np.random.randn(50), index=rng)

In [412]: resampled = data.resample('A', kind='timestamp', convention='end')

In [413]: resampled.index
Out[413]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2000-12-31 00:00:00, ..., 2004-12-31 00:00:00]
Length: 5, Freq: A-DEC, Timezone: None
```



# PACKAGE OVERVIEW

`pandas` consists of the following things

- A set of labeled array data structures, the primary of which are `Series/TimeSeries` and `DataFrame`
- Index objects enabling both simple axis indexing and multi-level / hierarchical axis indexing
- An integrated group by engine for aggregating and transforming data sets
- Date range generation (`date_range`) and custom date offsets enabling the implementation of customized frequencies
- Input/Output tools: loading tabular data from flat files (CSV, delimited, Excel 2003), and saving and loading `pandas` objects from the fast and efficient `PyTables/HDF5` format.
- Memory-efficient “sparse” versions of the standard data structures for storing data that is mostly missing or mostly constant (some fixed value)
- Moving window statistics (rolling mean, rolling standard deviation, etc.)
- Static and moving window linear and `panel regression`

## 4.1 Data structures at a glance

Dimensions	Name	Description
1	Series	1D labeled homogeneously-typed array
1	Time-Series	Series with index containing datetimes
2	DataFrame	General 2D labeled, size-mutable tabular structure with potentially heterogeneously-typed columns
3	Panel	General 3D labeled, also size-mutable array

### 4.1.1 Why more than 1 data structure?

The best way to think about the `pandas` data structures is as flexible containers for lower dimensional data. For example, `DataFrame` is a container for `Series`, and `Panel` is a container for `DataFrame` objects. We would like to be able to insert and remove objects from these containers in a dictionary-like fashion.

Also, we would like sensible default behaviors for the common API functions which take into account the typical orientation of time series and cross-sectional data sets. When using `ndarrays` to store 2- and 3-dimensional data, a burden is placed on the user to consider the orientation of the data set when writing functions; axes are considered more or less equivalent (except when C- or Fortran-contiguosness matters for performance). In `pandas`, the axes are

intended to lend more semantic meaning to the data; i.e., for a particular data set there is likely to be a “right” way to orient the data. The goal, then, is to reduce the amount of mental effort required to code up data transformations in downstream functions.

For example, with tabular data (DataFrame) it is more semantically helpful to think of the **index** (the rows) and the **columns** rather than axis 0 and axis 1. And iterating through the columns of the DataFrame thus results in more readable code:

```
for col in df.columns:
    series = df[col]
    # do something with series
```

## 4.2 Mutability and copying of data

All pandas data structures are value-mutable (the values they contain can be altered) but not always size-mutable. The length of a Series cannot be changed, but, for example, columns can be inserted into a DataFrame. However, the vast majority of methods produce new objects and leave the input data untouched. In general, though, we like to **favor immutability** where sensible.

## 4.3 Getting Support

The first stop for pandas issues and ideas is the [Github Issue Tracker](#). If you have a general question, pandas community experts can answer through [Stack Overflow](#).

Longer discussions occur on the [developer mailing list](#), and commercial support inquiries for Lambda Foundry should be sent to: [support@lambdafoundry.com](mailto:support@lambdafoundry.com)

## 4.4 Credits

pandas development began at [AQR Capital Management](#) in April 2008. It was open-sourced at the end of 2009. AQR continued to provide resources for development through the end of 2011, and continues to contribute bug reports today.

Since January 2012, [Lambda Foundry](#), has been providing development resources, as well as commercial support, training, and consulting for pandas.

pandas is only made possible by a group of people around the world like you who have contributed new code, bug reports, fixes, comments and ideas. A complete list can be found [on Github](#).

## 4.5 Development Team

pandas is a part of the PyData project. The PyData Development Team is a collection of developers focused on the improvement of Python’s data libraries. The core team that coordinates development can be found on [Github](#). If you’re interested in contributing, please visit the [project website](#).

## 4.6 License

=====  
PANDAS LICENSING TERMS  
=====

pandas is licensed under the BSD 3-Clause (also known as "BSD New" or "BSD Simplified"), as follows:

Copyright (c) 2011-2012, Lambda Foundry, Inc. and PyData Development Team  
All rights reserved.

Copyright (c) 2008-2011 AQR Capital Management, LLC  
All rights reserved.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

- \* Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- \* Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
- \* Neither the name of the copyright holder nor the names of any contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDER AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

About the Copyright Holders  
=====

AQR Capital Management began pandas development in 2008. Development was led by Wes McKinney. AQR released the source under this license in 2009. Wes is now an employee of Lambda Foundry, and remains the pandas project lead.

The PyData Development Team is the collection of developers of the PyData project. This includes all of the PyData sub-projects, including pandas. The core team that coordinates development on GitHub can be found here:  
<http://github.com/pydata>.

Full credits for pandas contributors can be found in the documentation.

Our Copyright Policy  
=====

PyData uses a shared copyright model. Each contributor maintains copyright over their contributions to PyData. However, it is important to note that these contributions are typically only changes to the repositories. Thus, the PyData source code, in its entirety, is not the copyright of any single person or institution. Instead, it is the collective copyright of the entire PyData Development Team. If individual contributors want to maintain a record of what changes/contributions they have specific copyright on, they should indicate their copyright in the commit message of the change when they commit the change to one of the PyData repositories.

With this in mind, the following banner should be used in any source code file to indicate the copyright and license terms:

```
#-----  
# Copyright (c) 2012, PyData Development Team  
# All rights reserved.  
#  
# Distributed under the terms of the BSD Simplified License.  
#  
# The full license is in the LICENSE file, distributed with this software.  
#-----
```

# INTRO TO DATA STRUCTURES

We'll start with a quick, non-comprehensive overview of the fundamental data structures in pandas to get you started. The fundamental behavior about data types, indexing, and axis labeling / alignment apply across all of the objects. To get started, import numpy and load pandas into your namespace:

```
In [261]: import numpy as np

# will use a lot in examples
In [262]: randn = np.random.randn

In [263]: from pandas import *
```

Here is a basic tenet to keep in mind: **data alignment is intrinsic**. Link between labels and data will not be broken unless done so explicitly by you.

We'll give a brief intro to the data structures, then consider all of the broad categories of functionality and methods in separate sections.

When using pandas, we recommend the following import convention:

```
import pandas as pd
```

## 5.1 Series

`Series` is a one-dimensional labeled array (technically a subclass of `ndarray`) capable of holding any data type (integers, strings, floating point numbers, Python objects, etc.). The axis labels are collectively referred to as the **index**. The basic method to create a `Series` is to call:

```
>>> s = Series(data, index=index)
```

Here, `data` can be many different things:

- a Python dict
- an `ndarray`
- a scalar value (like 5)

The passed **index** is a list of axis labels. Thus, this separates into a few cases depending on what **data** is:

### From `ndarray`

If `data` is an `ndarray`, **index** must be the same length as **data**. If no index is passed, one will be created having values `[0, ..., len(data) - 1]`.

```
In [264]: s = Series(randn(5), index=['a', 'b', 'c', 'd', 'e'])
```

```
In [265]: s
```

```
Out[265]:  
a    0.664  
b   -0.487  
c   -0.504  
d    0.307  
e    1.570
```

```
In [266]: s.index
```

```
Out[266]: Index([a, b, c, d, e], dtype=object)
```

```
In [267]: Series(randn(5))
```

```
Out[267]:  
0   -0.431  
1   -0.705  
2    0.555  
3    0.939  
4    0.722
```

---

**Note:** Starting in v0.8.0, pandas supports non-unique index values. In previous version, if the index values are not unique an exception will **not** be raised immediately, but attempting any operation involving the index will later result in an exception. In other words, the Index object containing the labels “lazily” checks whether the values are unique. The reason for being lazy is nearly all performance-based (there are many instances in computations, like parts of GroupBy, where the index is not used).

---

### From dict

If data is a dict, if **index** is passed the values in data corresponding to the labels in the index will be pulled out. Otherwise, an index will be constructed from the sorted keys of the dict, if possible.

```
In [268]: d = {'a' : 0., 'b' : 1., 'c' : 2.}
```

```
In [269]: Series(d)
```

```
Out[269]:  
a    0  
b    1  
c    2
```

```
In [270]: Series(d, index=['b', 'c', 'd', 'a'])
```

```
Out[270]:  
b    1  
c    2  
d   NaN  
a    0
```

---

**Note:** NaN (not a number) is the standard missing data marker used in pandas

---

**From scalar value** If data is a scalar value, an index must be provided. The value will be repeated to match the length of **index**

```
In [271]: Series(5., index=['a', 'b', 'c', 'd', 'e'])
```

```
Out[271]:  
a    5  
b    5
```

```
c    5
d    5
e    5
```

### 5.1.1 Series is ndarray-like

As a subclass of ndarray, Series is a valid argument to most NumPy functions and behaves similarly to a NumPy array. However, things like slicing also slice the index.

```
In [272]: s[0]
Out[272]: 0.66444516201494186
```

```
In [273]: s[:3]
Out[273]:
a    0.664
b   -0.487
c   -0.504
```

```
In [274]: s[s > s.median()]
Out[274]:
a    0.664
e    1.570
```

```
In [275]: s[[4, 3, 1]]
Out[275]:
e    1.570
d    0.307
b   -0.487
```

```
In [276]: np.exp(s)
Out[276]:
a    1.943
b    0.614
c    0.604
d    1.359
e    4.807
```

We will address array-based indexing in a separate [section](#).

### 5.1.2 Series is dict-like

A Series is alike a fixed-size dict in that you can get and set values by index label:

```
In [277]: s['a']
Out[277]: 0.66444516201494186
```

```
In [278]: s['e'] = 12.
```

```
In [279]: s
Out[279]:
a    0.664
b   -0.487
c   -0.504
d    0.307
e   12.000
```

```
In [280]: 'e' in s
Out[280]: True
```

```
In [281]: 'f' in s
Out[281]: False
```

If a label is not contained, an exception

```
>>> s['f']
KeyError: 'f'

>>> s.get('f')
nan
```

### 5.1.3 Vectorized operations and label alignment with Series

When doing data analysis, as with raw NumPy arrays looping through Series value-by-value is usually not necessary. Series can be also be passed into most NumPy methods expecting an ndarray.

```
In [282]: s + s
Out[282]:
a      1.329
b     -0.974
c     -1.009
d      0.613
e     24.000
```

```
In [283]: s * 2
Out[283]:
a      1.329
b     -0.974
c     -1.009
d      0.613
e     24.000
```

```
In [284]: np.exp(s)
Out[284]:
a      1.943
b      0.614
c      0.604
d      1.359
e    162754.791
```

A key difference between Series and ndarray is that operations between Series automatically align the data based on label. Thus, you can write computations without giving consideration to whether the Series involved have the same labels.

```
In [285]: s[1:] + s[:-1]
Out[285]:
a      NaN
b     -0.974
c     -1.009
d      0.613
e      NaN
```

The result of an operation between unaligned Series will have the **union** of the indexes involved. If a label is not found in one Series or the other, the result will be marked as missing (NaN). Being able to write code without doing any explicit data alignment grants immense freedom and flexibility in interactive data analysis and research. The integrated



data alignment features of the pandas data structures set pandas apart from the majority of related tools for working with labeled data.

---

**Note:** In general, we chose to make the default result of operations between differently indexed objects yield the **union** of the indexes in order to avoid loss of information. Having an index label, though the data is missing, is typically important information as part of a computation. You of course have the option of dropping labels with missing data via the **dropna** function.

---

### 5.1.4 Name attribute

Series can also have a name attribute:

```
In [286]: s = Series(np.random.randn(5), name='something')
```

```
In [287]: s
```

```
Out[287]:
```

```
0    0.015
```

```
1    1.987
```

```
2   -0.259
```

```
3    0.111
```

```
4    1.012
```

```
Name: something
```

```
In [288]: s.name
```

```
Out[288]: 'something'
```

The Series name will be assigned automatically in many cases, in particular when taking 1D slices of DataFrame as you will see below.

## 5.2 DataFrame

**DataFrame** is a 2-dimensional labeled data structure with columns of potentially different types. You can think of it like a spreadsheet or SQL table, or a dict of Series objects. It is generally the most commonly used pandas object. Like Series, DataFrame accepts many different kinds of input:

- Dict of 1D ndarrays, lists, dicts, or Series
- 2-D numpy.ndarray
- [Structured or record ndarray](#)
- A Series
- Another DataFrame

Along with the data, you can optionally pass **index** (row labels) and **columns** (column labels) arguments. If you pass an index and / or columns, you are guaranteeing the index and / or columns of the resulting DataFrame. Thus, a dict of Series plus a specific index will discard all data not matching up to the passed index.

If axis labels are not passed, they will be constructed from the input data based on common sense rules.

### 5.2.1 From dict of Series or dicts

The result **index** will be the **union** of the indexes of the various Series. If there are any nested dicts, these will be first converted to Series. If no columns are passed, the columns will be the sorted list of dict keys.

```
In [289]: d = {'one' : Series([1., 2., 3.], index=['a', 'b', 'c']),
.....:        'two' : Series([1., 2., 3., 4.], index=['a', 'b', 'c', 'd'])}
.....:
```

```
In [290]: df = DataFrame(d)
```

```
In [291]: df
```

```
Out[291]:
```

	one	two
a	1	1
b	2	2
c	3	3
d	NaN	4

```
In [292]: DataFrame(d, index=['d', 'b', 'a'])
```

```
Out[292]:
```

	one	two
d	NaN	4
b	2	2
a	1	1

```
In [293]: DataFrame(d, index=['d', 'b', 'a'], columns=['two', 'three'])
```

```
Out[293]:
```

	two	three
d	4	NaN
b	2	NaN
a	1	NaN

The row and column labels can be accessed respectively by accessing the **index** and **columns** attributes:

---

**Note:** When a particular set of columns is passed along with a dict of data, the passed columns override the keys in the dict.

---

```
In [294]: df.index
```

```
Out[294]: Index([a, b, c, d], dtype=object)
```

```
In [295]: df.columns
```

```
Out[295]: Index([one, two], dtype=object)
```

## 5.2.2 From dict of ndarrays / lists

The ndarrays must all be the same length. If an index is passed, it must clearly also be the same length as the arrays. If no index is passed, the result will be `range(n)`, where `n` is the array length.

```
In [296]: d = {'one' : [1., 2., 3., 4.],
.....:        'two' : [4., 3., 2., 1.]}
.....:
```

```
In [297]: DataFrame(d)
```

```
Out[297]:
```

	one	two
0	1	4
1	2	3
2	3	2
3	4	1

```
In [298]: DataFrame(d, index=['a', 'b', 'c', 'd'])
```

```
Out[298]:
```

	one	two
a	1	4
b	2	3
c	3	2
d	4	1

### 5.2.3 From structured or record array

This case is handled identically to a dict of arrays.

```
In [299]: data = np.zeros((2,), dtype=[('A', 'i4'), ('B', 'f4'), ('C', 'a10')])
```

```
In [300]: data[:] = [(1, 2., 'Hello'), (2, 3., "World")]
```

```
In [301]: DataFrame(data)
```

```
Out[301]:
```

	A	B	C
0	1	2	Hello
1	2	3	World

```
In [302]: DataFrame(data, index=['first', 'second'])
```

```
Out[302]:
```

	A	B	C
first	1	2	Hello
second	2	3	World

```
In [303]: DataFrame(data, columns=['C', 'A', 'B'])
```

```
Out[303]:
```

	C	A	B
0	Hello	1	2
1	World	2	3

---

**Note:** DataFrame is not intended to work exactly like a 2-dimensional NumPy ndarray.

---

### 5.2.4 From a list of dicts

```
In [304]: data2 = [{'a': 1, 'b': 2}, {'a': 5, 'b': 10, 'c': 20}]
```

```
In [305]: DataFrame(data2)
```

```
Out[305]:
```

	a	b	c
0	1	2	NaN
1	5	10	20

```
In [306]: DataFrame(data2, index=['first', 'second'])
```

```
Out[306]:
```

	a	b	c
first	1	2	NaN
second	5	10	20

```
In [307]: DataFrame(data2, columns=['a', 'b'])
```

```
Out[307]:
```

```
   a    b
0  1    2
1  5   10
```

### 5.2.5 From a Series

The result will be a DataFrame with the same index as the input Series, and with one column whose name is the original name of the Series (only if no other column name provided).

#### Missing Data

Much more will be said on this topic in the *Missing data* section. To construct a DataFrame with missing data, use `np.nan` for those values which are missing. Alternatively, you may pass a `numpy.MaskedArray` as the data argument to the DataFrame constructor, and its masked entries will be considered missing.

### 5.2.6 Alternate Constructors

#### DataFrame.from\_dict

`DataFrame.from_dict` takes a dict of dicts or a dict of array-like sequences and returns a DataFrame. It operates like the DataFrame constructor except for the `orient` parameter which is `'columns'` by default, but which can be set to `'index'` in order to use the dict keys as row labels. **DataFrame.from\_records**

`DataFrame.from_records` takes a list of tuples or an ndarray with structured dtype. Works analogously to the normal DataFrame constructor, except that index maybe be a specific field of the structured dtype to use as the index. For example:

```
In [308]: data
Out[308]:
array([(1, 2.0, 'Hello'), (2, 3.0, 'World')],
      dtype=[('A', '<i4'), ('B', '<f4'), ('C', '|S10')])
```

```
In [309]: DataFrame.from_records(data, index='C')
Out[309]:
   A  B
Hello 1  2
World 2  3
```

#### DataFrame.from\_items

`DataFrame.from_items` works analogously to the form of the dict constructor that takes a sequence of (key, value) pairs, where the keys are column (or row, in the case of `orient='index'`) names, and the value are the column values (or row values). This can be useful for constructing a DataFrame with the columns in a particular order without having to pass an explicit list of columns:

```
In [310]: DataFrame.from_items([('A', [1, 2, 3]), ('B', [4, 5, 6])])
Out[310]:
   A  B
0  1  4
1  2  5
2  3  6
```

If you pass `orient='index'`, the keys will be the row labels. But in this case you must also pass the desired column names:

```
In [311]: DataFrame.from_items([(‘A’, [1, 2, 3]), (‘B’, [4, 5, 6])],
.....:                        orient=‘index’, columns=[‘one’, ‘two’, ‘three’])
.....:
Out[311]:
```

	one	two	three
A	1	2	3
B	4	5	6

### 5.2.7 Column selection, addition, deletion

You can treat a DataFrame semantically like a dict of like-indexed Series objects. Getting, setting, and deleting columns works with the same syntax as the analogous dict operations:

```
In [312]: df[‘one’]
Out[312]:
```

a	1
b	2
c	3
d	NaN

Name: one

```
In [313]: df[‘three’] = df[‘one’] * df[‘two’]

In [314]: df[‘flag’] = df[‘one’] > 2
```

```
In [315]: df
Out[315]:
```

	one	two	three	flag
a	1	1	1	False
b	2	2	4	False
c	3	3	9	True
d	NaN	4	NaN	False

Columns can be deleted or popped like with a dict:

```
In [316]: del df[‘two’]

In [317]: three = df.pop(‘three’)

In [318]: df
Out[318]:
```

	one	flag
a	1	False
b	2	False
c	3	True
d	NaN	False

When inserting a scalar value, it will naturally be propagated to fill the column:

```
In [319]: df[‘foo’] = ‘bar’

In [320]: df
Out[320]:
```

	one	flag	foo
a	1	False	bar
b	2	False	bar
c	3	True	bar
d	NaN	False	bar

When inserting a Series that does not have the same index as the DataFrame, it will be conformed to the DataFrame's index:

```
In [321]: df['one_trunc'] = df['one'][:2]
```

```
In [322]: df
```

```
Out[322]:
```

	one	flag	foo	one_trunc
a	1	False	bar	1
b	2	False	bar	2
c	3	True	bar	NaN
d	NaN	False	bar	NaN

You can insert raw ndarrays but their length must match the length of the DataFrame's index.

By default, columns get inserted at the end. The `insert` function is available to insert at a particular location in the columns:

```
In [323]: df.insert(1, 'bar', df['one'])
```

```
In [324]: df
```

```
Out[324]:
```

	one	bar	flag	foo	one_trunc
a	1	1	False	bar	1
b	2	2	False	bar	2
c	3	3	True	bar	NaN
d	NaN	NaN	False	bar	NaN

## 5.2.8 Indexing / Selection

The basics of indexing are as follows:

Operation	Syntax	Result
Select column	<code>df[col]</code>	Series
Select row by label	<code>df.xs(label)</code> or <code>df.ix[label]</code>	Series
Select row by location (int)	<code>df.ix[loc]</code>	Series
Slice rows	<code>df[5:10]</code>	DataFrame
Select rows by boolean vector	<code>df[bool_vec]</code>	DataFrame

Row selection, for example, returns a Series whose index is the columns of the DataFrame:

```
In [325]: df.xs('b')
```

```
Out[325]:
```

one	2
bar	2
flag	False
foo	bar
one_trunc	2

Name: b

```
In [326]: df.ix[2]
```

```
Out[326]:
```

one	3
bar	3
flag	True
foo	bar
one_trunc	NaN

Name: c

Note if a DataFrame contains columns of multiple dtypes, the dtype of the row will be chosen to accommodate all of the data types (dtype=object is the most general).

For a more exhaustive treatment of more sophisticated label-based indexing and slicing, see the [section on indexing](#). We will address the fundamentals of reindexing / conforming to new sets of labels in the [section on reindexing](#).

## 5.2.9 Data alignment and arithmetic

Data alignment between DataFrame objects automatically align on **both the columns and the index (row labels)**. Again, the resulting object will have the union of the column and row labels.

```
In [327]: df = DataFrame(randn(10, 4), columns=['A', 'B', 'C', 'D'])
```

```
In [328]: df2 = DataFrame(randn(7, 3), columns=['A', 'B', 'C'])
```

```
In [329]: df + df2
```

```
Out [329]:
```

	A	B	C	D
0	2.752	-0.429	0.702	NaN
1	0.067	-3.397	1.775	NaN
2	-0.499	-1.138	-1.277	NaN
3	0.731	0.988	0.505	NaN
4	-0.538	-1.828	-1.974	NaN
5	-0.100	-2.885	1.676	NaN
6	1.405	-1.078	0.320	NaN
7	NaN	NaN	NaN	NaN
8	NaN	NaN	NaN	NaN
9	NaN	NaN	NaN	NaN

When doing an operation between DataFrame and Series, the default behavior is to align the Series **index** on the DataFrame **columns**, thus **broadcasting** row-wise. For example:

```
In [330]: df - df.ix[0]
```

```
Out [330]:
```

	A	B	C	D
0	0.000	0.000	0.000	0.000
1	-0.586	-3.234	-1.119	-2.876
2	-1.136	-3.727	-2.066	-1.500
3	-0.799	-0.803	-1.860	-2.038
4	-1.632	-2.216	-2.384	-1.367
5	-0.485	-2.545	-1.044	-1.190
6	-0.168	0.395	-1.333	-1.060
7	-2.367	-0.799	-3.479	-2.653
8	-2.953	-1.565	-2.251	-2.978
9	-1.615	-1.712	-2.521	-2.106

In the special case of working with time series data, if the Series is a TimeSeries (which it will be automatically if the index contains datetime objects), and the DataFrame index also contains dates, the broadcasting will be column-wise:

```
In [331]: index = date_range('1/1/2000', periods=8)
```

```
In [332]: df = DataFrame(randn(8, 3), index=index,
.....:                  columns=['A', 'B', 'C'])
.....:
```

```
In [333]: df
```

```
Out [333]:
```

	A	B	C
--	---	---	---

```
2000-01-01 -1.209 -1.257 -0.500
2000-01-02  0.430 -0.242 -0.724
2000-01-03  1.257 -0.871 -0.544
2000-01-04 -0.766 -0.219  0.663
2000-01-05 -1.566  1.780 -2.139
2000-01-06 -0.593 -1.059  0.119
2000-01-07 -0.123  1.306 -0.129
2000-01-08 -0.389  0.143 -1.715
```

```
In [334]: type(df['A'])
Out[334]: pandas.core.series.TimeSeries
```

```
In [335]: df - df['A']
Out[335]:
```

	A	B	C
2000-01-01	0	-0.048	0.709
2000-01-02	0	-0.672	-1.154
2000-01-03	0	-2.128	-1.801
2000-01-04	0	0.547	1.429
2000-01-05	0	3.346	-0.572
2000-01-06	0	-0.466	0.711
2000-01-07	0	1.429	-0.006
2000-01-08	0	0.532	-1.326

Technical purity aside, this case is so common in practice that supporting the special case is preferable to the alternative of forcing the user to transpose and do column-based alignment like so:

```
In [336]: (df.T - df['A']).T
Out[336]:
```

	A	B	C
2000-01-01	0	-0.048	0.709
2000-01-02	0	-0.672	-1.154
2000-01-03	0	-2.128	-1.801
2000-01-04	0	0.547	1.429
2000-01-05	0	3.346	-0.572
2000-01-06	0	-0.466	0.711
2000-01-07	0	1.429	-0.006
2000-01-08	0	0.532	-1.326

For explicit control over the matching and broadcasting behavior, see the section on *flexible binary operations*.

Operations with scalars are just as you would expect:

```
In [337]: df * 5 + 2
Out[337]:
```

	A	B	C
2000-01-01	-4.043	-4.285	-0.499
2000-01-02	4.149	0.789	-1.619
2000-01-03	8.286	-2.355	-0.719
2000-01-04	-1.830	0.907	5.314
2000-01-05	-5.831	10.900	-8.693
2000-01-06	-0.963	-3.296	2.594
2000-01-07	1.385	8.528	1.354
2000-01-08	0.054	2.716	-6.575

```
In [338]: 1 / df
Out[338]:
```

	A	B	C
2000-01-01	-0.827	-0.796	-2.001



```

2000-01-02    2.327 -4.129 -1.382
2000-01-03    0.795 -1.148 -1.839
2000-01-04   -1.305 -4.574  1.509
2000-01-05   -0.638  0.562 -0.468
2000-01-06   -1.687 -0.944  8.416
2000-01-07   -8.128  0.766 -7.746
2000-01-08   -2.570  6.983 -0.583

```

```
In [339]: df ** 4
```

```
Out [339]:
```

```

          A          B          C
2000-01-01  2.133    2.497    0.062
2000-01-02  0.034    0.003    0.275
2000-01-03  2.499    0.576    0.087
2000-01-04  0.344    0.002    0.193
2000-01-05  6.018   10.038   20.918
2000-01-06  0.123    1.258    0.000
2000-01-07  0.000    2.906    0.000
2000-01-08  0.023    0.000    8.652

```

Boolean operators work as well:

```
In [340]: df1 = DataFrame({'a' : [1, 0, 1], 'b' : [0, 1, 1] }, dtype=bool)
```

```
In [341]: df2 = DataFrame({'a' : [0, 1, 1], 'b' : [1, 1, 0] }, dtype=bool)
```

```
In [342]: df1 & df2
```

```
Out [342]:
```

```

      a      b
0  False False
1  False  True
2   True False

```

```
In [343]: df1 | df2
```

```
Out [343]:
```

```

      a      b
0   True  True
1   True  True
2   True  True

```

```
In [344]: df1 ^ df2
```

```
Out [344]:
```

```

      a      b
0   True  True
1   True False
2  False  True

```

```
In [345]: ~df1
```

```
Out [345]:
```

```

      a      b
0  False  True
1   True False
2  False False

```

## 5.2.10 Transposing

To transpose, access the `T` attribute (also the `transpose` function), similar to an `ndarray`:

```
# only show the first 5 rows
In [346]: df[:5].T
Out[346]:
```

	2000-01-01	2000-01-02	2000-01-03	2000-01-04	2000-01-05
A	-1.209	0.430	1.257	-0.766	-1.566
B	-1.257	-0.242	-0.871	-0.219	1.780
C	-0.500	-0.724	-0.544	0.663	-2.139

### 5.2.11 DataFrame interoperability with NumPy functions

Elementwise NumPy ufuncs (log, exp, sqrt, ...) and various other NumPy functions can be used with no issues on DataFrame, assuming the data within are numeric:

```
In [347]: np.exp(df)
Out[347]:
```

	A	B	C
2000-01-01	0.299	0.285	0.607
2000-01-02	1.537	0.785	0.485
2000-01-03	3.516	0.419	0.581
2000-01-04	0.465	0.804	1.940
2000-01-05	0.209	5.930	0.118
2000-01-06	0.553	0.347	1.126
2000-01-07	0.884	3.690	0.879
2000-01-08	0.678	1.154	0.180

```
In [348]: np.asarray(df)
Out[348]:
```

```
array([[ -1.2085,  -1.257 ,  -0.4997],
       [  0.4298,  -0.2422,  -0.7238],
       [  1.2573,  -0.871 ,  -0.5437],
       [-0.7661,  -0.2186,   0.6628],
       [-1.5663,   1.78   , -2.1386],
       [-0.5927,  -1.0591,   0.1188],
       [-0.123 ,   1.3056,  -0.1291],
       [-0.3892,   0.1432,  -1.715 ]])
```

The dot method on DataFrame implements matrix multiplication:

```
In [349]: df.T.dot(df)
Out[349]:
```

	A	B	C
A	6.784	-1.889	3.064
B	-1.889	8.460	-3.214
C	3.064	-3.214	9.055

Similarly, the dot method on Series implements dot product:

```
In [350]: s1 = Series(np.arange(5,10))

In [351]: s1.dot(s1)
Out[351]: 255
```

DataFrame is not intended to be a drop-in replacement for ndarray as its indexing semantics are quite different in places from a matrix.

## 5.2.12 Console display

For very large DataFrame objects, only a summary will be printed to the console (here I am reading a CSV version of the **baseball** dataset from the **plyr** R package):

```
In [352]: baseball = read_csv('data/baseball.csv')
```

```
In [353]: print baseball
<class 'pandas.core.frame.DataFrame'>
Int64Index: 100 entries, 88641 to 89534
Data columns:
id          100  non-null values
year        100  non-null values
stint       100  non-null values
team        100  non-null values
lg          100  non-null values
g           100  non-null values
ab          100  non-null values
r           100  non-null values
h           100  non-null values
X2b         100  non-null values
X3b         100  non-null values
hr          100  non-null values
rbi         100  non-null values
sb          100  non-null values
cs          100  non-null values
bb          100  non-null values
so          100  non-null values
ibb         100  non-null values
hbp         100  non-null values
sh          100  non-null values
sf          100  non-null values
gidp        100  non-null values
dtypes: float64(9), int64(10), object(3)
```

However, using `to_string` will return a string representation of the DataFrame in tabular form, though it won't always fit the console width:

```
In [354]: print baseball.ix[-20:, :12].to_string()
      id  year  stint team  lg   g  ab   r   h  X2b  X3b  hr
88641  womacto01  2006     2  CHN  NL   19  50    6  14    1    0    1
88643  schilcu01  2006     1  BOS  AL   31    2    0    1    0    0    0
88645  myersmi01  2006     1  NYA  AL   62    0    0    0    0    0    0
88649  helliri01  2006     1  MIL  NL   20    3    0    0    0    0    0
88650  johnsra05  2006     1  NYA  AL   33    6    0    1    0    0    0
88652  finlest01  2006     1  SFN  NL  139  426   66  105   21   12    6
88653  gonzalu01  2006     1  ARI  NL  153  586   93  159   52    2   15
88662  seleaa01  2006     1  LAN  NL   28   26    2    5    1    0    0
89177  francju01  2007     2  ATL  NL   15   40    1   10    3    0    0
89178  francju01  2007     1  NYN  NL   40   50    7   10    0    0    1
89330  zaungr01  2007     1  TOR  AL  110  331   43   80   24    1   10
89333  witasja01  2007     1  TBA  AL    3    0    0    0    0    0    0
89334  williwo02  2007     1  HOU  NL   33   59    3    6    0    0    1
89335  wickmbo01  2007     2  ARI  NL    8    0    0    0    0    0    0
89336  wickmbo01  2007     1  ATL  NL   47    0    0    0    0    0    0
89337  whitero02  2007     1  MIN  AL   38  109    8   19    4    0    4
89338  whiteri01  2007     1  HOU  NL   20    1    0    0    0    0    0
89339  wellsda01  2007     2  LAN  NL    7   15    2    4    1    0    0
89340  wellsda01  2007     1  SDN  NL   22   38    1    4    0    0    0
```

89341	weathda01	2007	1	CIN	NL	67	0	0	0	0	0
89343	walketo04	2007	1	OAK	AL	18	48	5	13	1	0
89345	wakefti01	2007	1	BOS	AL	1	2	0	0	0	0
89347	vizquom01	2007	1	SFN	NL	145	513	54	126	18	3
89348	villoro01	2007	1	NYA	AL	6	0	0	0	0	0
89352	valenjo03	2007	1	NYN	NL	51	166	18	40	11	1
89354	trachst01	2007	2	CHN	NL	4	7	0	1	0	0
89355	trachst01	2007	1	BAL	AL	3	5	0	0	0	0
89359	timlimi01	2007	1	BOS	AL	4	0	0	0	0	0
89360	thomeji01	2007	1	CHA	AL	130	432	79	119	19	0
89361	thomafr04	2007	1	TOR	AL	155	531	63	147	30	0
89363	tavarju01	2007	1	BOS	AL	2	4	0	1	0	0
89365	sweenma01	2007	2	LAN	NL	30	33	2	9	1	0
89366	sweenma01	2007	1	SFN	NL	76	90	18	23	8	0
89367	suppaje01	2007	1	MIL	NL	33	61	4	8	0	0
89368	stinnke01	2007	1	SLN	NL	26	82	7	13	3	0
89370	stantmi02	2007	1	CIN	NL	67	2	0	0	0	0
89371	stairma01	2007	1	TOR	AL	125	357	58	103	28	1
89372	sprinru01	2007	1	SLN	NL	72	1	0	0	0	0
89374	sosasa01	2007	1	TEX	AL	114	412	53	104	24	1
89375	smoltjo01	2007	1	ATL	NL	30	54	1	5	1	0
89378	sheffga01	2007	1	DET	AL	133	494	107	131	20	1
89381	seleaa01	2007	1	NYN	NL	31	4	0	0	0	0
89382	seaneru01	2007	1	LAN	NL	68	1	0	0	0	0
89383	schmija01	2007	1	LAN	NL	6	7	1	1	0	0
89384	schilcu01	2007	1	BOS	AL	1	2	0	1	0	0
89385	sandere02	2007	1	KCA	AL	24	73	12	23	7	0
89388	rogerke01	2007	1	DET	AL	1	2	0	0	0	0
89389	rodriiv01	2007	1	DET	AL	129	502	50	141	31	3
89396	ramirma02	2007	1	BOS	AL	133	483	84	143	33	1
89398	piazzmi01	2007	1	OAK	AL	83	309	33	85	17	1
89400	perezne01	2007	1	DET	AL	33	64	5	11	3	0
89402	parkch01	2007	1	NYN	NL	1	1	0	0	0	0
89406	oliveda02	2007	1	LAA	AL	5	0	0	0	0	0
89410	myersmi01	2007	1	NYA	AL	6	1	0	0	0	0
89411	mussimi01	2007	1	NYA	AL	2	2	0	0	0	0
89412	moyerja01	2007	1	PHI	NL	33	73	4	9	2	0
89420	mesajo01	2007	1	PHI	NL	38	0	0	0	0	0
89421	martipe02	2007	1	NYN	NL	5	9	1	1	1	0
89425	maddugr01	2007	1	SDN	NL	33	62	2	9	2	0
89426	mabryjo01	2007	1	COL	NL	28	34	4	4	1	0
89429	loftoke01	2007	2	CLE	AL	52	173	24	49	9	3
89430	loftoke01	2007	1	TEX	AL	84	317	62	96	16	3
89431	loaizes01	2007	1	LAN	NL	5	7	0	1	0	0
89438	kleskry01	2007	1	SFN	NL	116	362	51	94	27	3
89439	kentje01	2007	1	LAN	NL	136	494	78	149	36	1
89442	jonesto02	2007	1	DET	AL	5	0	0	0	0	0
89445	johnsra05	2007	1	ARI	NL	10	15	0	1	0	0
89450	hoffmtr01	2007	1	SDN	NL	60	0	0	0	0	0
89451	hernaro01	2007	2	LAN	NL	22	0	0	0	0	0
89452	hernaro01	2007	1	CLE	AL	2	0	0	0	0	0
89460	guarded01	2007	1	CIN	NL	15	0	0	0	0	0
89462	griffke02	2007	1	CIN	NL	144	528	78	146	24	1
89463	greensh01	2007	1	NYN	NL	130	446	62	130	30	1
89464	graffto01	2007	1	MIL	NL	86	231	34	55	8	0
89465	gordoto01	2007	1	PHI	NL	44	0	0	0	0	0
89466	gonzalu01	2007	1	LAN	NL	139	464	70	129	23	2
89467	gomezch02	2007	2	CLE	AL	19	53	4	15	2	0

89468	gomezch02	2007	1	BAL	AL	73	169	17	51	10	1	1
89469	glavito02	2007	1	NYN	NL	33	56	3	12	1	0	0
89473	floydc101	2007	1	CHN	NL	108	282	40	80	10	1	9
89474	finlest01	2007	1	COL	NL	43	94	9	17	3	0	1
89480	embreal01	2007	1	OAK	AL	4	0	0	0	0	0	0
89481	edmonji01	2007	1	SLN	NL	117	365	39	92	15	2	12
89482	easleda01	2007	1	NYN	NL	76	193	24	54	6	0	10
89489	delgaca01	2007	1	NYN	NL	139	538	71	139	30	0	24
89493	cormirh01	2007	1	CIN	NL	6	0	0	0	0	0	0
89494	coninje01	2007	2	NYN	NL	21	41	2	8	2	0	0
89495	coninje01	2007	1	CIN	NL	80	215	23	57	11	1	6
89497	clemero02	2007	1	NYA	AL	2	2	0	1	0	0	0
89498	claytro01	2007	2	BOS	AL	8	6	1	0	0	0	0
89499	claytro01	2007	1	TOR	AL	69	189	23	48	14	0	1
89501	cirilje01	2007	2	ARI	NL	28	40	6	8	4	0	0
89502	cirilje01	2007	1	MIN	AL	50	153	18	40	9	2	2
89521	bondsba01	2007	1	SFN	NL	126	340	75	94	14	0	28
89523	biggicr01	2007	1	HOU	NL	141	517	68	130	31	3	10
89525	benitar01	2007	2	FLO	NL	34	0	0	0	0	0	0
89526	benitar01	2007	1	SFN	NL	19	0	0	0	0	0	0
89530	ausmubr01	2007	1	HOU	NL	117	349	38	82	16	3	3
89533	aloumo01	2007	1	NYN	NL	87	328	51	112	19	1	13
89534	alomasa02	2007	1	NYN	NL	8	22	1	3	1	0	0

### 5.2.13 DataFrame column types

The four main types stored in pandas objects are float, int, boolean, and object. A convenient `dtypes` attribute return a Series with the data type of each column:

```
In [355]: baseball.dtypes
```

```
Out[355]:
id          object
year        int64
stint       int64
team        object
lg          object
g           int64
ab          int64
r           int64
h           int64
X2b         int64
X3b         int64
hr          int64
rbi         float64
sb          float64
cs          float64
bb          int64
so          float64
ibb         float64
hbp         float64
sh          float64
sf          float64
gidp        float64
```

The related method `get_dtype_counts` will return the number of columns of each type:

```
In [356]: baseball.get_dtype_counts()
Out[356]:
float64      9
int64        10
object        3
```

### 5.2.14 DataFrame column attribute access and IPython completion

If a DataFrame column label is a valid Python variable name, the column can be accessed like attributes:

```
In [357]: df = DataFrame({'foo1' : np.random.randn(5),
.....:                  'foo2' : np.random.randn(5)})
.....:
```

```
In [358]: df
Out[358]:
      foo1      foo2
0  0.759091 -0.648742
1 -0.050457  0.209870
2  0.959219 -0.325391
3 -0.817600 -1.978199
4 -0.200407 -0.211127
```

```
In [359]: df.foo1
Out[359]:
0    0.759091
1   -0.050457
2    0.959219
3   -0.817600
4   -0.200407
Name: foo1
```

The columns are also connected to the IPython completion mechanism so they can be tab-completed:

```
In [5]: df.fo<TAB>
df.foo1  df.foo2
```

## 5.3 Panel

Panel is a somewhat less-used, but still important container for 3-dimensional data. The term **panel data** is derived from econometrics and is partially responsible for the name pandas: pan(el)-da(ta)-s. The names for the 3 axes are intended to give some semantic meaning to describing operations involving panel data and, in particular, econometric analysis of panel data. However, for the strict purposes of slicing and dicing a collection of DataFrame objects, you may find the axis names slightly arbitrary:

- **items**: axis 0, each item corresponds to a DataFrame contained inside
- **major\_axis**: axis 1, it is the **index** (rows) of each of the DataFrames
- **minor\_axis**: axis 2, it is the **columns** of each of the DataFrames

Construction of Panels works about like you would expect:

### 5.3.1 From 3D ndarray with optional axis labels

```
In [360]: wp = Panel(randn(2, 5, 4), items=['Item1', 'Item2'],
.....:               major_axis=date_range('1/1/2000', periods=5),
.....:               minor_axis=['A', 'B', 'C', 'D'])
.....:
```

```
In [361]: wp
Out[361]:
<class 'pandas.core.panel.Panel'>
Dimensions: 2 (items) x 5 (major) x 4 (minor)
Items: Item1 to Item2
Major axis: 2000-01-01 00:00:00 to 2000-01-05 00:00:00
Minor axis: A to D
```

### 5.3.2 From dict of DataFrame objects

```
In [362]: data = {'Item1' : DataFrame(randn(4, 3)),
.....:           'Item2' : DataFrame(randn(4, 2))}
.....:
```

```
In [363]: Panel(data)
Out[363]:
<class 'pandas.core.panel.Panel'>
Dimensions: 2 (items) x 4 (major) x 3 (minor)
Items: Item1 to Item2
Major axis: 0 to 3
Minor axis: 0 to 2
```

Note that the values in the dict need only be **convertible to DataFrame**. Thus, they can be any of the other valid inputs to DataFrame as per above.

One helpful factory method is `Panel.from_dict`, which takes a dictionary of DataFrames as above, and the following named parameters:

Parameter	Default	Description
<code>intersect</code>	<code>False</code>	drops elements whose indices do not align
<code>orient</code>	<code>items</code>	use <code>minor</code> to use DataFrames' columns as panel items

For example, compare to the construction above:

```
In [364]: Panel.from_dict(data, orient='minor')
Out[364]:
<class 'pandas.core.panel.Panel'>
Dimensions: 3 (items) x 4 (major) x 2 (minor)
Items: 0 to 2
Major axis: 0 to 3
Minor axis: Item1 to Item2
```

Orient is especially useful for mixed-type DataFrames. If you pass a dict of DataFrame objects with mixed-type columns, all of the data will get upcasted to `dtype=object` unless you pass `orient='minor'`:

```
In [365]: df = DataFrame({'a': ['foo', 'bar', 'baz'],
.....:                  'b': np.random.randn(3)})
.....:
```

```
In [366]: df
Out[366]:
```

```
      a      b
0  foo  0.080597
1  bar -0.000185
2  baz -0.264704
```

```
In [367]: data = {'item1': df, 'item2': df}
```

```
In [368]: panel = Panel.from_dict(data, orient='minor')
```

```
In [369]: panel['a']
```

```
Out[369]:
      item1 item2
0     foo     foo
1     bar     bar
2     baz     baz
```

```
In [370]: panel['b']
```

```
Out[370]:
      item1      item2
0  0.080597  0.080597
1 -0.000185 -0.000185
2 -0.264704 -0.264704
```

```
In [371]: panel['b'].dtypes
```

```
Out[371]:
item1      float64
item2      float64
```

---

**Note:** Unfortunately Panel, being less commonly used than Series and DataFrame, has been slightly neglected feature-wise. A number of methods and options available in DataFrame are not available in Panel. This will get worked on, of course, in future releases. And faster if you join me in working on the codebase.

---

### 5.3.3 From DataFrame using `to_panel` method

This method was introduced in v0.7 to replace `LongPanel.to_long`, and converts a DataFrame with a two-level index to a Panel.

```
In [372]: midx = MultiIndex(levels=[['one', 'two'], ['x', 'y']], labels=[[1,1,0,0],[1,0,1,0]])
```

```
In [373]: df = DataFrame({'A' : [1, 2, 3, 4], 'B': [5, 6, 7, 8]}, index=midx)
```

```
In [374]: df.to_panel()
```

```
Out[374]:
<class 'pandas.core.panel.Panel'>
Dimensions: 2 (items) x 2 (major) x 2 (minor)
Items: A to B
Major axis: one to two
Minor axis: x to y
```

### 5.3.4 Item selection / addition / deletion

Similar to DataFrame functioning as a dict of Series, Panel is like a dict of DataFrames:



```
In [375]: wp['Item1']
```

```
Out[375]:
```

	A	B	C	D
2000-01-01	-0.519332	-1.765523	-0.966196	-0.890524
2000-01-02	-1.314597	-1.458515	-0.919663	-0.699091
2000-01-03	1.357258	-0.098278	-0.987183	-1.362030
2000-01-04	-1.309989	-1.153000	0.606382	-0.681101
2000-01-05	-0.289724	-0.996632	-1.407699	1.014104

```
In [376]: wp['Item3'] = wp['Item1'] / wp['Item2']
```

The API for insertion and deletion is the same as for DataFrame. And as with DataFrame, if the item is a valid python identifier, you can access it as an attribute and tab-complete it in IPython.

### 5.3.5 Transposing

A Panel can be rearranged using its `transpose` method (which does not make a copy by default unless the data are heterogeneous):

```
In [377]: wp.transpose(2, 0, 1)
```

```
Out[377]:
```

```
<class 'pandas.core.panel.Panel'>
Dimensions: 4 (items) x 3 (major) x 5 (minor)
Items: A to D
Major axis: Item1 to Item3
Minor axis: 2000-01-01 00:00:00 to 2000-01-05 00:00:00
```

### 5.3.6 Indexing / Selection

Operation	Syntax	Result
Select item	<code>wp[item]</code>	DataFrame
Get slice at major_axis label	<code>wp.major_xs(val)</code>	DataFrame
Get slice at minor_axis label	<code>wp.minor_xs(val)</code>	DataFrame

For example, using the earlier example data, we could do:

```
In [378]: wp['Item1']
```

```
Out[378]:
```

	A	B	C	D
2000-01-01	-0.519332	-1.765523	-0.966196	-0.890524
2000-01-02	-1.314597	-1.458515	-0.919663	-0.699091
2000-01-03	1.357258	-0.098278	-0.987183	-1.362030
2000-01-04	-1.309989	-1.153000	0.606382	-0.681101
2000-01-05	-0.289724	-0.996632	-1.407699	1.014104

```
In [379]: wp.major_xs(wp.major_axis[2])
```

```
Out[379]:
```

	Item1	Item2	Item3
A	1.357258	-0.177665	-7.639427
B	-0.098278	0.490838	-0.200224
C	-0.987183	-1.360102	0.725815
D	-1.362030	1.592456	-0.855302

```
In [380]: wp.minor_axis
```

```
Out[380]: Index([A, B, C, D], dtype=object)
```

```
In [381]: wp.minor_xs('C')
Out[381]:
```

	Item1	Item2	Item3
2000-01-01	-0.966196	0.071823	-13.452418
2000-01-02	-0.919663	0.214910	-4.279288
2000-01-03	-0.987183	-1.360102	0.725815
2000-01-04	0.606382	-1.890591	-0.320737
2000-01-05	-1.407699	-0.151652	9.282439

### 5.3.7 Conversion to DataFrame

A Panel can be represented in 2D form as a hierarchically indexed DataFrame. See the section *hierarchical indexing* for more on this. To convert a Panel to a DataFrame, use the `to_frame` method:

```
In [382]: panel = Panel(np.random.randn(3, 5, 4), items=['one', 'two', 'three'],
.....:                  major_axis=date_range('1/1/2000', periods=5),
.....:                  minor_axis=['a', 'b', 'c', 'd'])
.....:
```

```
In [383]: panel.to_frame()
Out[383]:
```

		one	two	three
2000-01-01	a	-0.566820	0.597468	0.716659
	b	-1.643966	-0.491240	-0.919717
	c	1.471262	1.281674	-0.024595
	d	0.677634	-0.099685	0.068997
2000-01-02	a	-0.485743	-1.823043	0.601797
	b	-0.342272	-0.779213	0.866615
	c	-1.042291	-0.949327	0.092911
	d	-0.611457	0.768043	-2.606892
2000-01-03	a	-0.141224	-0.054860	0.309303
	b	0.007220	-1.493561	-0.548401
	c	-0.516147	0.106004	-2.044772
	d	0.446161	-0.903513	-1.666264
2000-01-04	a	0.483368	-0.719875	-1.439775
	b	0.186405	0.301945	1.326361
	c	-1.439567	1.112546	0.221680
	d	-0.503782	-0.542770	1.840992
2000-01-05	a	0.890769	-2.695540	1.165150
	b	-0.777798	0.431284	-1.420521
	c	-0.552820	-0.431092	1.616679
	d	-1.428744	1.666631	-1.030912

---

# ESSENTIAL BASIC FUNCTIONALITY

Here we discuss a lot of the essential functionality common to the pandas data structures. Here's how to create some of the objects used in the examples from the previous section:

```
In [1]: index = date_range('1/1/2000', periods=8)

In [2]: s = Series(randn(5), index=['a', 'b', 'c', 'd', 'e'])

In [3]: df = DataFrame(randn(8, 3), index=index,
...:                   columns=['A', 'B', 'C'])
...:
...:

In [4]: wp = Panel(randn(2, 5, 4), items=['Item1', 'Item2'],
...:                major_axis=date_range('1/1/2000', periods=5),
...:                minor_axis=['A', 'B', 'C', 'D'])
...:
...:
```

## 6.1 Head and Tail

To view a small sample of a Series or DataFrame object, use the `head` and `tail` methods. The default number of elements to display is five, but you may pass a custom number.

```
In [5]: long_series = Series(randn(1000))

In [6]: long_series.head()
Out[6]:
0      0.951142
1      0.262153
2     -1.472473
3     -0.935746
4      0.951939

In [7]: long_series.tail(3)
Out[7]:
997      1.108191
998      0.025259
999      0.523356
```

## 6.2 Attributes and the raw ndarray(s)

pandas objects have a number of attributes enabling you to access the metadata

- **shape**: gives the axis dimensions of the object, consistent with ndarray
- Axis labels
  - **Series**: *index* (only axis)
  - **DataFrame**: *index* (rows) and *columns*
  - **Panel**: *items*, *major\_axis*, and *minor\_axis*

Note, these attributes can be safely assigned to!

```
In [8]: df[:2]
```

```
Out[8]:
```

	A	B	C
2000-01-01	-1.242120	-0.053063	1.143213
2000-01-02	0.557477	-1.148352	1.590601

```
In [9]: df.columns = [x.lower() for x in df.columns]
```

```
In [10]: df
```

```
Out[10]:
```

	a	b	c
2000-01-01	-1.242120	-0.053063	1.143213
2000-01-02	0.557477	-1.148352	1.590601
2000-01-03	-0.451055	-0.084402	1.204146
2000-01-04	-0.917944	0.077257	1.347842
2000-01-05	0.384912	-1.095539	0.361094
2000-01-06	-0.979042	1.176231	-0.261979
2000-01-07	0.070772	0.974250	1.010776
2000-01-08	-0.028392	1.989038	-0.566515

To get the actual data inside a data structure, one need only access the **values** property:

```
In [11]: s.values
```

```
Out[11]: array([ 0.7697,  1.762 ,  0.4605,  0.9344,  1.2188])
```

```
In [12]: df.values
```

```
Out[12]:
```

```
array([[ -1.2421,  -0.0531,   1.1432],
       [  0.5575,  -1.1484,   1.5906],
       [ -0.4511,  -0.0844,   1.2041],
       [ -0.9179,   0.0773,   1.3478],
       [  0.3849,  -1.0955,   0.3611],
       [ -0.979 ,   1.1762,  -0.262 ],
       [  0.0708,   0.9742,   1.0108],
       [ -0.0284,   1.989 ,  -0.5665]])
```

```
In [13]: wp.values
```

```
Out[13]:
```

```
array([[ -0.6244,  -0.3029,   1.7733,   2.6527],
       [ -2.0124,  -0.1523,  -0.2526,   0.9716],
       [ -1.1932,   0.5684,  -0.0966,   0.3392],
       [  1.4384,   0.0409,  -1.9785,   0.4886],
       [  0.4925,   1.2875,   0.3125,  -0.269 ],
       [ -1.1927,  -0.8258,  -1.4743,  -0.1465],
       [ -0.5977,  -1.1085,   0.673 ,  -0.4669],
```

```
[ 0.0869,  0.1202,  1.3991, -1.807 ],
[ 0.86   ,  0.0182,  1.2952, -0.8072],
[ 0.9639, -1.0762,  0.3464,  0.4477]]])
```

If a DataFrame or Panel contains homogeneously-typed data, the ndarray can actually be modified in-place, and the changes will be reflected in the data structure. For heterogeneous data (e.g. some of the DataFrame's columns are not all the same dtype), this will not be the case. The values attribute itself, unlike the axis labels, cannot be assigned to.

---

**Note:** When working with heterogeneous data, the dtype of the resulting ndarray will be chosen to accommodate all of the data involved. For example, if strings are involved, the result will be of object dtype. If there are only floats and integers, the resulting array will be of float dtype.

---

## 6.3 Flexible binary operations

With binary operations between pandas data structures, there are two key points of interest:

- Broadcasting behavior between higher- (e.g. DataFrame) and lower-dimensional (e.g. Series) objects.
- Missing data in computations

We will demonstrate how to manage these issues independently, though they can be handled simultaneously.

### 6.3.1 Matching / broadcasting behavior

DataFrame has the methods **add**, **sub**, **mul**, **div** and related functions **radd**, **rsub**, ... for carrying out binary operations. For broadcasting behavior, Series input is of primary interest. Using these functions, you can use to either match on the *index* or *columns* via the **axis** keyword:

```
In [14]: df
Out[14]:
```

	one	three	two
a	0.588637	NaN	-0.478462
b	-1.167858	0.085539	0.719219
c	0.524038	0.910929	1.129957
d	NaN	-1.274663	0.739634

```
In [15]: row = df.ix[1]

In [16]: column = df['two']

In [17]: df.sub(row, axis='columns')
Out[17]:
```

	one	three	two
a	1.756495	NaN	-1.197681
b	0.000000	0.000000	0.000000
c	1.691896	0.825390	0.410738
d	NaN	-1.360202	0.020415

```
In [18]: df.sub(row, axis=1)
Out[18]:
```

	one	three	two
a	1.756495	NaN	-1.197681
b	0.000000	0.000000	0.000000
c	1.691896	0.825390	0.410738

```
d      NaN -1.360202  0.020415
```

```
In [19]: df.sub(column, axis='index')
```

```
Out[19]:
```

	one	three	two
a	1.067099	NaN	0
b	-1.887077	-0.633680	0
c	-0.605919	-0.219028	0
d	NaN	-2.014297	0

```
In [20]: df.sub(column, axis=0)
```

```
Out[20]:
```

	one	three	two
a	1.067099	NaN	0
b	-1.887077	-0.633680	0
c	-0.605919	-0.219028	0
d	NaN	-2.014297	0

With Panel, describing the matching behavior is a bit more difficult, so the arithmetic methods instead (and perhaps confusingly?) give you the option to specify the *broadcast axis*. For example, suppose we wished to demean the data over a particular axis. This can be accomplished by taking the mean over an axis and broadcasting over the same axis:

```
In [21]: major_mean = wp.mean(axis='major')
```

```
In [22]: major_mean
```

```
Out[22]:
```

	Item1	Item2
A	-0.379820	0.024078
B	0.288297	-0.574406
C	-0.048374	0.447876
D	0.836640	-0.555944

```
In [23]: wp.sub(major_mean, axis='major')
```

```
Out[23]:
```

```
<class 'pandas.core.panel.Panel'>
Dimensions: 2 (items) x 5 (major) x 4 (minor)
Items: Item1 to Item2
Major axis: 2000-01-01 00:00:00 to 2000-01-05 00:00:00
Minor axis: A to D
```

And similarly for `axis="items"` and `axis="minor"`.

---

**Note:** I could be convinced to make the **axis** argument in the DataFrame methods match the broadcasting behavior of Panel. Though it would require a transition period so users can change their code...

---

## 6.3.2 Missing data / operations with fill values

In Series and DataFrame (though not yet in Panel), the arithmetic functions have the option of inputting a *fill\_value*, namely a value to substitute when at most one of the values at a location are missing. For example, when adding two DataFrame objects, you may wish to treat NaN as 0 unless both DataFrames are missing that value, in which case the result will be NaN (you can later replace NaN with some other value using `fillna` if you wish).

```
In [24]: df
```

```
Out[24]:
```

	one	three	two
a	0.588637	NaN	-0.478462

```
b -1.167858  0.085539  0.719219
c  0.524038  0.910929  1.129957
d           NaN -1.274663  0.739634
```

```
In [25]: df2
```

```
Out [25]:
```

	one	three	two
a	0.588637	1.000000	-0.478462
b	-1.167858	0.085539	0.719219
c	0.524038	0.910929	1.129957
d	NaN	-1.274663	0.739634

```
In [26]: df + df2
```

```
Out [26]:
```

	one	three	two
a	1.177275	NaN	-0.956923
b	-2.335716	0.171079	1.438438
c	1.048076	1.821859	2.259915
d	NaN	-2.549326	1.479268

```
In [27]: df.add(df2, fill_value=0)
```

```
Out [27]:
```

	one	three	two
a	1.177275	1.000000	-0.956923
b	-2.335716	0.171079	1.438438
c	1.048076	1.821859	2.259915
d	NaN	-2.549326	1.479268

### 6.3.3 Flexible Comparisons

Starting in v0.8, pandas introduced binary comparison methods `eq`, `ne`, `lt`, `gt`, `le`, and `ge` to `Series` and `DataFrame` whose behavior is analogous to the binary arithmetic operations described above:

```
In [28]: df.gt(df2)
```

```
Out [28]:
```

	one	three	two
a	False	False	False
b	False	False	False
c	False	False	False
d	False	False	False

```
In [29]: df2.ne(df)
```

```
Out [29]:
```

	one	three	two
a	False	True	False
b	False	False	False
c	False	False	False
d	True	False	False

### 6.3.4 Combining overlapping data sets

A problem occasionally arising is the combination of two similar data sets where values in one are preferred over the other. An example would be two data series representing a particular economic indicator where one is considered to be of “higher quality”. However, the lower quality series might extend further back in history or have more complete data coverage. As such, we would like to combine two `DataFrame` objects where missing values in one `DataFrame`

are conditionally filled with like-labeled values from the other DataFrame. The function implementing this operation is `combine_first`, which we illustrate:

```
In [30]: df1 = DataFrame({'A' : [1., np.nan, 3., 5., np.nan],
.....:                  'B' : [np.nan, 2., 3., np.nan, 6.]})
.....:

In [31]: df2 = DataFrame({'A' : [5., 2., 4., np.nan, 3., 7.],
.....:                  'B' : [np.nan, np.nan, 3., 4., 6., 8.]})
.....:
```

```
In [32]: df1
```

```
Out[32]:
```

	A	B
0	1	NaN
1	NaN	2
2	3	3
3	5	NaN
4	NaN	6

```
In [33]: df2
```

```
Out[33]:
```

	A	B
0	5	NaN
1	2	NaN
2	4	3
3	NaN	4
4	3	6
5	7	8

```
In [34]: df1.combine_first(df2)
```

```
Out[34]:
```

	A	B
0	1	NaN
1	2	2
2	3	3
3	5	4
4	3	6
5	7	8

### 6.3.5 General DataFrame Combine

The `combine_first` method above calls the more general DataFrame method `combine`. This method takes another DataFrame and a combiner function, aligns the input DataFrame and then passes the combiner function pairs of Series (ie, columns whose names are the same).

So, for instance, to reproduce `combine_first` as above:

```
In [35]: combiner = lambda x, y: np.where(isnull(x), y, x)
```

```
In [36]: df1.combine(df2, combiner)
```

```
Out[36]:
```

	A	B
0	1	NaN
1	2	2
2	3	3
3	5	4



```
4  3    6
5  7    8
```

## 6.4 Descriptive statistics

A large number of methods for computing descriptive statistics and other related operations on *Series*, *DataFrame*, and *Panel*. Most of these are aggregations (hence producing a lower-dimensional result) like **sum**, **mean**, and **quantile**, but some of them, like **cumsum** and **cumprod**, produce an object of the same size. Generally speaking, these methods take an **axis** argument, just like *ndarray*.{*sum*, *std*, ...}, but the axis can be specified by name or integer:

- **Series**: no axis argument needed
- **DataFrame**: “index” (axis=0, default), “columns” (axis=1)
- **Panel**: “items” (axis=0), “major” (axis=1, default), “minor” (axis=2)

For example:

```
In [37]: df
Out[37]:
```

	one	three	two
a	0.588637	NaN	-0.478462
b	-1.167858	0.085539	0.719219
c	0.524038	0.910929	1.129957
d	NaN	-1.274663	0.739634

```
In [38]: df.mean(0)
Out[38]:
```

one	-0.018394
three	-0.092731
two	0.527587

```
In [39]: df.mean(1)
Out[39]:
```

a	0.055088
b	-0.121033
c	0.854975
d	-0.267514

All such methods have a `skipna` option signaling whether to exclude missing data (`True` by default):

```
In [40]: df.sum(0, skipna=False)
Out[40]:
```

one	NaN
three	NaN
two	2.110349

```
In [41]: df.sum(axis=1, skipna=True)
Out[41]:
```

a	0.110176
b	-0.363099
c	2.564925
d	-0.535029

Combined with the broadcasting / arithmetic behavior, one can describe various statistical procedures, like standardization (rendering data zero mean and standard deviation 1), very concisely:

```
In [42]: ts_stand = (df - df.mean()) / df.std()
```

```
In [43]: ts_stand.std()
```

```
Out[43]:  
one      1  
three    1  
two      1
```

```
In [44]: xs_stand = df.sub(df.mean(1), axis=0).div(df.std(1), axis=0)
```

```
In [45]: xs_stand.std(1)
```

```
Out[45]:  
a      1  
b      1  
c      1  
d      1
```

Note that methods like **cumsum** and **cumprod** preserve the location of NA values:

```
In [46]: df.cumsum()
```

```
Out[46]:  
      one      three      two  
a  0.588637      NaN -0.478462  
b -0.579221  0.085539  0.240758  
c -0.055182  0.996469  1.370715  
d      NaN -0.278194  2.110349
```

Here is a quick reference summary table of common functions. Each also takes an optional `level` parameter which applies only if the object has a *hierarchical index*.

Function	Description
count	Number of non-null observations
sum	Sum of values
mean	Mean of values
mad	Mean absolute deviation
median	Arithmetic median of values
min	Minimum
max	Maximum
abs	Absolute Value
prod	Product of values
std	Unbiased standard deviation
var	Unbiased variance
skew	Unbiased skewness (3rd moment)
kurt	Unbiased kurtosis (4th moment)
quantile	Sample quantile (value at %)
cumsum	Cumulative sum
cumprod	Cumulative product
cummax	Cumulative maximum
cummin	Cumulative minimum

Note that by chance some NumPy methods, like `mean`, `std`, and `sum`, will exclude NAs on Series input by default:

```
In [47]: np.mean(df['one'])
```

```
Out[47]: -0.018394106418168483
```

```
In [48]: np.mean(df['one'].values)
```

```
Out[48]: nan
```

Series also has a method `nunique` which will return the number of unique non-null values:

```
In [49]: series = Series(randn(500))
```

```
In [50]: series[20:500] = np.nan
```

```
In [51]: series[10:20] = 5
```

```
In [52]: series.nunique()
```

```
Out[52]: 11
```

### 6.4.1 Summarizing data: describe

There is a convenient `describe` function which computes a variety of summary statistics about a Series or the columns of a DataFrame (excluding NAs of course):

```
In [53]: series = Series(randn(1000))
```

```
In [54]: series[::2] = np.nan
```

```
In [55]: series.describe()
```

```
Out[55]:
count    500.000000
mean      0.023591
std       1.101907
min      -3.043228
25%      -0.722430
50%       0.040272
75%       0.742909
max       3.307501
```

```
In [56]: frame = DataFrame(randn(1000, 5), columns=['a', 'b', 'c', 'd', 'e'])
```

```
In [57]: frame.ix[::2] = np.nan
```

```
In [58]: frame.describe()
```

```
Out[58]:
```

	a	b	c	d	e
count	500.000000	500.000000	500.000000	500.000000	500.000000
mean	0.046406	-0.002076	0.037099	-0.008660	0.002823
std	0.958654	0.981637	0.981859	1.061440	1.047552
min	-2.692483	-2.795762	-2.600736	-3.204610	-3.233743
25%	-0.645598	-0.682636	-0.603275	-0.702792	-0.688680
50%	0.067905	-0.042679	0.056899	0.029435	0.034855
75%	0.664905	0.600600	0.744798	0.726872	0.712177
max	3.065924	3.439255	3.442551	3.285388	3.093917

For a non-numerical Series object, `describe` will give a simple summary of the number of unique values and most frequently occurring values:

```
In [59]: s = Series(['a', 'a', 'b', 'b', 'a', 'a', np.nan, 'c', 'd', 'a'])
```

```
In [60]: s.describe()
```

```
Out[60]:
count    9
unique    4
top      a
freq     5
```

There also is a utility function, `value_range` which takes a `DataFrame` and returns a series with the minimum/maximum values in the `DataFrame`.

## 6.4.2 Index of Min/Max Values

The `idxmin` and `idxmax` functions on `Series` and `DataFrame` compute the index labels with the minimum and maximum corresponding values:

```
In [61]: s1 = Series(randn(5))
```

```
In [62]: s1
```

```
Out[62]:
0    -1.227467
1    -0.684019
2     0.952409
3     1.343835
4    -0.072712
```

```
In [63]: s1.idxmin(), s1.idxmax()
```

```
Out[63]: (0, 3)
```

```
In [64]: df1 = DataFrame(randn(5,3), columns=['A','B','C'])
```

```
In [65]: df1
```

```
Out[65]:
      A         B         C
0 -1.372590 -1.367802 -0.028855
1  1.374353  0.603209  0.823161
2 -0.604128 -1.231778  1.938518
3 -0.217089  1.082010 -0.930506
4  1.031590  0.631951 -0.789942
```

```
In [66]: df1.idxmin(axis=0)
```

```
Out[66]:
```

```
A    0
B    0
C    3
```

```
In [67]: df1.idxmax(axis=1)
```

```
Out[67]:
```

```
0    C
1    A
2    C
3    B
4    A
```

When there are multiple rows (or columns) matching the minimum or maximum value, `idxmin` and `idxmax` return the first matching index:

```
In [68]: df3 = DataFrame([2, 1, 1, 3, np.nan], columns=['A'], index=list('edcba'))
```

```
In [69]: df3
```

```
Out[69]:
```

```
      A
e     2
d     1
c     1
b     3
```

a NaN

```
In [70]: df3['A'].idxmin()
Out[70]: 'd'
```

### 6.4.3 Value counts (histogramming)

The `value_counts` Series method and top-level function computes a histogram of a 1D array of values. It can also be used as a function on regular arrays:

```
In [71]: data = np.random.randint(0, 7, size=50)
```

```
In [72]: data
Out[72]:
array([0, 3, 2, 0, 0, 2, 3, 2, 0, 3, 0, 5, 1, 4, 3, 4, 2, 3, 1, 1, 2, 1, 5,
        6, 4, 5, 0, 2, 4, 3, 2, 1, 1, 3, 2, 0, 4, 3, 1, 6, 6, 1, 0, 1, 2, 0,
        4, 0, 1, 0])
```

```
In [73]: s = Series(data)
```

```
In [74]: s.value_counts()
```

```
Out[74]:
0    11
1    10
2     9
3     8
4     6
6     3
5     3
```

```
In [75]: value_counts(data)
```

```
Out[75]:
0    11
1    10
2     9
3     8
4     6
6     3
5     3
```

### 6.4.4 Discretization and quantiling

Continuous values can be discretized using the `cut` (bins based on values) and `qcut` (bins based on sample quantiles) functions:

```
In [76]: arr = np.random.randn(20)
```

```
In [77]: factor = cut(arr, 4)
```

```
In [78]: factor
```

```
Out[78]:
Categorical:
array([(1.015, 2.397], (-0.367, 1.015], (-0.367, 1.015], (-1.749, -0.367],
       (1.015, 2.397], (-0.367, 1.015], (1.015, 2.397], (-0.367, 1.015],
       (-0.367, 1.015], (-3.136, -1.749], (-0.367, 1.015], (-0.367, 1.015],
```

```
(-0.367, 1.015], (-0.367, 1.015], (1.015, 2.397], (-3.136, -1.749],
(-0.367, 1.015], (-3.136, -1.749], (1.015, 2.397], (-0.367, 1.015]], dtype=object)
Levels (4): Index([(-3.136, -1.749], (-1.749, -0.367], (-0.367, 1.015],
(1.015, 2.397]], dtype=object)
```

```
In [79]: factor = cut(arr, [-5, -1, 0, 1, 5])
```

```
In [80]: factor
```

```
Out[80]:
```

```
Categorical:
```

```
array([(1, 5], (0, 1], (0, 1], (-5, -1], (1, 5], (0, 1], (1, 5], (-1, 0],
      (0, 1], (-5, -1], (0, 1], (0, 1], (0, 1], (0, 1], (1, 5], (-5, -1],
      (0, 1], (-5, -1], (1, 5], (0, 1]], dtype=object)
```

```
Levels (4): Index([(-5, -1], (-1, 0], (0, 1], (1, 5]], dtype=object)
```

`qcut` computes sample quantiles. For example, we could slice up some normally distributed data into equal-size quantiles like so:

```
In [81]: arr = np.random.randn(30)
```

```
In [82]: factor = qcut(arr, [0, .25, .5, .75, 1])
```

```
In [83]: factor
```

```
Out[83]:
```

```
Categorical:
```

```
array([(0.068, 0.716], (0.716, 2.252], [-1.222, -0.585], [-1.222, -0.585],
      (0.716, 2.252], (0.716, 2.252], (0.716, 2.252], [-1.222, -0.585],
      (-0.585, 0.068], (0.068, 0.716], (0.068, 0.716], (-0.585, 0.068],
      [-1.222, -0.585], (0.716, 2.252], [-1.222, -0.585], (0.716, 2.252],
      (0.716, 2.252], [-1.222, -0.585], (-0.585, 0.068], (0.068, 0.716],
      (0.068, 0.716], (-0.585, 0.068], (-0.585, 0.068], (0.068, 0.716],
      (-0.585, 0.068], [-1.222, -0.585], (0.716, 2.252], (0.068, 0.716],
      [-1.222, -0.585], (-0.585, 0.068]], dtype=object)
```

```
Levels (4): Index([[-1.222, -0.585], (-0.585, 0.068], (0.068, 0.716],
(0.716, 2.252]], dtype=object)
```

```
In [84]: value_counts(factor)
```

```
Out[84]:
```

```
[-1.222, -0.585]    8
(0.716, 2.252]      8
(0.068, 0.716]      7
(-0.585, 0.068]     7
```

## 6.5 Function application

Arbitrary functions can be applied along the axes of a `DataFrame` or `Panel` using the `apply` method, which, like the descriptive statistics methods, take an optional `axis` argument:

```
In [85]: df.apply(np.mean)
```

```
Out[85]:
```

```
one      -0.018394
three    -0.092731
two       0.527587
```

```
In [86]: df.apply(np.mean, axis=1)
```

```
Out[86]:
```

```
a    0.055088
b   -0.121033
c    0.854975
d   -0.267514
```

```
In [87]: df.apply(lambda x: x.max() - x.min())
```

```
Out[87]:
one      1.756495
three    2.185592
two      1.608419
```

```
In [88]: df.apply(np.cumsum)
```

```
Out[88]:
      one      three      two
a  0.588637      NaN -0.478462
b -0.579221  0.085539  0.240758
c -0.055182  0.996469  1.370715
d      NaN -0.278194  2.110349
```

```
In [89]: df.apply(np.exp)
```

```
Out[89]:
      one      three      two
a  1.801532      NaN  0.619736
b  0.311032  1.089305  2.052830
c  1.688834  2.486633  3.095255
d      NaN  0.279525  2.095169
```

Depending on the return type of the function passed to `apply`, the result will either be of lower dimension or the same dimension.

`apply` combined with some cleverness can be used to answer many questions about a data set. For example, suppose we wanted to extract the date where the maximum value for each column occurred:

```
In [90]: tsdf = DataFrame(randn(1000, 3), columns=['A', 'B', 'C'],
.....:                      index=date_range('1/1/2000', periods=1000))
.....:
```

```
In [91]: tsdf.apply(lambda x: x.index[x.dropna().argmax()])
```

```
Out[91]:
A    2001-04-05 00:00:00
B    2002-06-26 00:00:00
C    2000-06-14 00:00:00
```

You may also pass additional arguments and keyword arguments to the `apply` method. For instance, consider the following function you would like to apply:

```
def subtract_and_divide(x, sub, divide=1):
    return (x - sub) / divide
```

You may then apply this function as follows:

```
df.apply(subtract_and_divide, args=(5,), divide=3)
```

Another useful feature is the ability to pass Series methods to carry out some Series operation on each column or row:

```
In [92]: tsdf
```

```
Out[92]:
      A      B      C
2000-01-01 -1.752694 -0.235040  0.534535
2000-01-02  1.535579 -0.733470 -0.344265
```

```
2000-01-03    1.601153    2.675714    0.895098
2000-01-04         NaN         NaN         NaN
2000-01-05         NaN         NaN         NaN
2000-01-06         NaN         NaN         NaN
2000-01-07         NaN         NaN         NaN
2000-01-08    0.000612    0.315425   -2.228711
2000-01-09   -0.543118    1.379495    1.544588
2000-01-10    0.213461   -0.561242    0.838548
```

```
In [93]: tsdf.apply(Series.interpolate)
```

```
Out [93]:
```

```
          A          B          C
2000-01-01 -1.752694 -0.235040  0.534535
2000-01-02  1.535579 -0.733470 -0.344265
2000-01-03  1.601153  2.675714  0.895098
2000-01-04  1.281045  2.203656  0.270337
2000-01-05  0.960937  1.731598 -0.354425
2000-01-06  0.640828  1.259540 -0.979187
2000-01-07  0.320720  0.787483 -1.603949
2000-01-08  0.000612  0.315425 -2.228711
2000-01-09 -0.543118  1.379495  1.544588
2000-01-10  0.213461 -0.561242  0.838548
```

Finally, `apply` takes an argument `raw` which is `False` by default, which converts each row or column into a `Series` before applying the function. When set to `True`, the passed function will instead receive an `ndarray` object, which has positive performance implications if you do not need the indexing functionality.

#### See Also:

The section on [GroupBy](#) demonstrates related, flexible functionality for grouping by some criterion, applying, and combining the results into a `Series`, `DataFrame`, etc.

## 6.5.1 Applying elementwise Python functions

Since not all functions can be vectorized (accept NumPy arrays and return another array or value), the methods `applymap` on `DataFrame` and analogously `map` on `Series` accept any Python function taking a single value and returning a single value. For example:

```
In [94]: f = lambda x: len(str(x))
```

```
In [95]: df['one'].map(f)
```

```
Out [95]:
```

```
a    14
b    14
c    14
d     3
Name: one
```

```
In [96]: df.applymap(f)
```

```
Out [96]:
```

```
   one  three  two
a   14     3   15
b   14    15   14
c   14    14   13
d    3    13   14
```

`Series.map` has an additional feature which is that it can be used to easily “link” or “map” values defined by a secondary series. This is closely related to [merging/joining functionality](#):



```
In [97]: s = Series(['six', 'seven', 'six', 'seven', 'six'],
.....:             index=['a', 'b', 'c', 'd', 'e'])
.....:
```

```
In [98]: t = Series({'six' : 6., 'seven' : 7.})
```

```
In [99]: s
```

```
Out[99]:
```

```
a      six
b     seven
c      six
d     seven
e      six
```

```
In [100]: s.map(t)
```

```
Out[100]:
```

```
a      6
b      7
c      6
d      7
e      6
```

## 6.6 Reindexing and altering labels

`reindex` is the fundamental data alignment method in pandas. It is used to implement nearly all other features relying on label-alignment functionality. To *reindex* means to conform the data to match a given set of labels along a particular axis. This accomplishes several things:

- Reorders the existing data to match a new set of labels
- Inserts missing value (NA) markers in label locations where no data for that label existed
- If specified, **fill** data for missing labels using logic (highly relevant to working with time series data)

Here is a simple example:

```
In [101]: s = Series(randn(5), index=['a', 'b', 'c', 'd', 'e'])
```

```
In [102]: s
```

```
Out[102]:
```

```
a      0.272664
b     -2.072136
c     -0.059607
d     -1.718696
e      0.128977
```

```
In [103]: s.reindex(['e', 'b', 'f', 'd'])
```

```
Out[103]:
```

```
e      0.128977
b     -2.072136
f           NaN
d     -1.718696
```

Here, the `f` label was not contained in the Series and hence appears as `NaN` in the result.

With a DataFrame, you can simultaneously reindex the index and columns:

```
In [104]: df
Out[104]:
```

	one	three	two
a	0.588637	NaN	-0.478462
b	-1.167858	0.085539	0.719219
c	0.524038	0.910929	1.129957
d	NaN	-1.274663	0.739634

```
In [105]: df.reindex(index=['c', 'f', 'b'], columns=['three', 'two', 'one'])
Out[105]:
```

	three	two	one
c	0.910929	1.129957	0.524038
f	NaN	NaN	NaN
b	0.085539	0.719219	-1.167858

For convenience, you may utilize the `reindex_axis` method, which takes the labels and a keyword `axis` parameter.

Note that the `Index` objects containing the actual axis labels can be **shared** between objects. So if we have a `Series` and a `DataFrame`, the following can be done:

```
In [106]: rs = s.reindex(df.index)

In [107]: rs
Out[107]:
```

a	0.272664
b	-2.072136
c	-0.059607
d	-1.718696

```
In [108]: rs.index is df.index
Out[108]: True
```

This means that the reindexed `Series`'s index is the same Python object as the `DataFrame`'s index.

#### See Also:

*Advanced indexing* is an even more concise way of doing reindexing.

---

**Note:** When writing performance-sensitive code, there is a good reason to spend some time becoming a reindexing ninja: **many operations are faster on pre-aligned data**. Adding two unaligned `DataFrames` internally triggers a reindexing step. For exploratory analysis you will hardly notice the difference (because `reindex` has been heavily optimized), but when CPU cycles matter sprinkling a few explicit `reindex` calls here and there can have an impact.

---

## 6.6.1 Reindexing to align with another object

You may wish to take an object and reindex its axes to be labeled the same as another object. While the syntax for this is straightforward albeit verbose, it is a common enough operation that the `reindex_like` method is available to make this simpler:

```
In [109]: df
Out[109]:
```

	one	three	two
a	0.588637	NaN	-0.478462
b	-1.167858	0.085539	0.719219
c	0.524038	0.910929	1.129957
d	NaN	-1.274663	0.739634

```
In [110]: df2
Out[110]:
```

	one	two
a	0.607032	-0.935367
b	-1.149464	0.262314
c	0.542432	0.673052

```
In [111]: df.reindex_like(df2)
Out[111]:
```

	one	two
a	0.588637	-0.478462
b	-1.167858	0.719219
c	0.524038	1.129957

## 6.6.2 Reindexing with `reindex_axis`

## 6.6.3 Aligning objects with each other with `align`

The `align` method is the fastest way to simultaneously align two objects. It supports a `join` argument (related to *joining and merging*):

- `join='outer'`: take the union of the indexes
- `join='left'`: use the calling object's index
- `join='right'`: use the passed object's index
- `join='inner'`: intersect the indexes

It returns a tuple with both of the reindexed Series:

```
In [112]: s = Series(randn(5), index=['a', 'b', 'c', 'd', 'e'])

In [113]: s1 = s[:4]

In [114]: s2 = s[1:]

In [115]: s1.align(s2)
Out[115]:
```

(a	0.684584
b	-0.612873
c	0.673347
d	2.283157
e	NaN,
a	NaN
b	-0.612873
c	0.673347
d	2.283157
e	1.165933)

```
In [116]: s1.align(s2, join='inner')
Out[116]:
```

(b	-0.612873
c	0.673347
d	2.283157,
b	-0.612873
c	0.673347
d	2.283157)

```
In [117]: s1.align(s2, join='left')
Out[117]:
(a      0.684584
 b     -0.612873
 c      0.673347
 d      2.283157,
 a           NaN
 b     -0.612873
 c      0.673347
 d      2.283157)
```

For DataFrames, the join method will be applied to both the index and the columns by default:

```
In [118]: df.align(df2, join='inner')
Out[118]:
(      one      two
a  0.588637 -0.478462
b -1.167858  0.719219
c  0.524038  1.129957,
      one      two
a  0.607032 -0.935367
b -1.149464  0.262314
c  0.542432  0.673052)
```

You can also pass an axis option to only align on the specified axis:

```
In [119]: df.align(df2, join='inner', axis=0)
Out[119]:
(      one      three      two
a  0.588637      NaN -0.478462
b -1.167858  0.085539  0.719219
c  0.524038  0.910929  1.129957,
      one      two
a  0.607032 -0.935367
b -1.149464  0.262314
c  0.542432  0.673052)
```

If you pass a Series to DataFrame.align, you can choose to align both objects either on the DataFrame's index or columns using the axis argument:

```
In [120]: df.align(df2.ix[0], axis=1)
Out[120]:
(      one      three      two
a  0.588637      NaN -0.478462
b -1.167858  0.085539  0.719219
c  0.524038  0.910929  1.129957
d           NaN -1.274663  0.739634,
  one      0.607032
three      NaN
two      -0.935367
Name: a)
```

## 6.6.4 Filling while reindexing

reindex takes an optional parameter method which is a filling method chosen from the following table:

Method	Action
pad / ffill	Fill values forward
bfill / backfill	Fill values backward

Other fill methods could be added, of course, but these are the two most commonly used for time series data. In a way they only make sense for time series or otherwise ordered data, but you may have an application on non-time series data where this sort of “interpolation” logic is the correct thing to do. More sophisticated interpolation of missing values would be an obvious extension.

We illustrate these fill methods on a simple TimeSeries:

```
In [121]: rng = date_range('1/3/2000', periods=8)
```

```
In [122]: ts = Series(randn(8), index=rng)
```

```
In [123]: ts2 = ts[[0, 3, 6]]
```

```
In [124]: ts
```

```
Out [124]:
2000-01-03    -1.732762
2000-01-04    -1.169055
2000-01-05     0.809870
2000-01-06     2.094194
2000-01-07     0.722281
2000-01-08     1.796709
2000-01-09     0.716624
2000-01-10     0.464436
Freq: D
```

```
In [125]: ts2
```

```
Out [125]:
2000-01-03    -1.732762
2000-01-06     2.094194
2000-01-09     0.716624
```

```
In [126]: ts2.reindex(ts.index)
```

```
Out [126]:
2000-01-03    -1.732762
2000-01-04         NaN
2000-01-05         NaN
2000-01-06     2.094194
2000-01-07         NaN
2000-01-08         NaN
2000-01-09     0.716624
2000-01-10         NaN
Freq: D
```

```
In [127]: ts2.reindex(ts.index, method='ffill')
```

```
Out [127]:
2000-01-03    -1.732762
2000-01-04    -1.732762
2000-01-05    -1.732762
2000-01-06     2.094194
2000-01-07     2.094194
2000-01-08     2.094194
2000-01-09     0.716624
2000-01-10     0.716624
Freq: D
```

```
In [128]: ts2.reindex(ts.index, method='bfill')
Out[128]:
2000-01-03    -1.732762
2000-01-04     2.094194
2000-01-05     2.094194
2000-01-06     2.094194
2000-01-07     0.716624
2000-01-08     0.716624
2000-01-09     0.716624
2000-01-10         NaN
Freq: D
```

Note the same result could have been achieved using *fillna*:

```
In [129]: ts2.reindex(ts.index).fillna(method='ffill')
Out[129]:
2000-01-03    -1.732762
2000-01-04    -1.732762
2000-01-05    -1.732762
2000-01-06     2.094194
2000-01-07     2.094194
2000-01-08     2.094194
2000-01-09     0.716624
2000-01-10     0.716624
Freq: D
```

Note these methods generally assume that the indexes are **sorted**. They may be modified in the future to be a bit more flexible but as time series data is ordered most of the time anyway, this has not been a major priority.

## 6.6.5 Dropping labels from an axis

A method closely related to `reindex` is the `drop` function. It removes a set of labels from an axis:

```
In [130]: df
Out[130]:
      one      three      two
a  0.588637      NaN -0.478462
b -1.167858  0.085539  0.719219
c  0.524038  0.910929  1.129957
d      NaN -1.274663  0.739634

In [131]: df.drop(['a', 'd'], axis=0)
Out[131]:
      one      three      two
b -1.167858  0.085539  0.719219
c  0.524038  0.910929  1.129957

In [132]: df.drop(['one'], axis=1)
Out[132]:
      three      two
a      NaN -0.478462
b  0.085539  0.719219
c  0.910929  1.129957
d -1.274663  0.739634
```

Note that the following also works, but is a bit less obvious / clean:

```
In [133]: df.reindex(df.index - ['a', 'd'])
Out[133]:
```

	one	three	two
b	-1.167858	0.085539	0.719219
c	0.524038	0.910929	1.129957

### 6.6.6 Renaming / mapping labels

The `rename` method allows you to relabel an axis based on some mapping (a dict or Series) or an arbitrary function.

```
In [134]: s
Out[134]:
```

a	0.684584
b	-0.612873
c	0.673347
d	2.283157
e	1.165933

```
In [135]: s.rename(str.upper)
Out[135]:
```

A	0.684584
B	-0.612873
C	0.673347
D	2.283157
E	1.165933

If you pass a function, it must return a value when called with any of the labels (and must produce a set of unique values). But if you pass a dict or Series, it need only contain a subset of the labels as keys:

```
In [136]: df.rename(columns={'one' : 'foo', 'two' : 'bar'},
.....:               index={'a' : 'apple', 'b' : 'banana', 'd' : 'durian'})
Out[136]:
```

	foo	three	bar
apple	0.588637	NaN	-0.478462
banana	-1.167858	0.085539	0.719219
c	0.524038	0.910929	1.129957
durian	NaN	-1.274663	0.739634

The `rename` method also provides an `inplace` named parameter that is by default `False` and copies the underlying data. Pass `inplace=True` to rename the data in place. The `Panel` class has a related `rename_axis` class which can rename any of its three axes.

## 6.7 Iteration

Because `Series` is array-like, basic iteration produces the values. Other data structures follow the dict-like convention of iterating over the “keys” of the objects. In short:

- **Series:** values
- **DataFrame:** column labels
- **Panel:** item labels

Thus, for example:

```
In [137]: for col in df:
.....:     print col
.....:
one
three
two
```

### 6.7.1 iteritems

Consistent with the dict-like interface, **iteritems** iterates through key-value pairs:

- **Series**: (index, scalar value) pairs
- **DataFrame**: (column, Series) pairs
- **Panel**: (item, DataFrame) pairs

For example:

```
In [138]: for item, frame in wp.iteritems():
.....:     print item
.....:     print frame
.....:
```

Item1

	A	B	C	D
2000-01-01	-0.624390	-0.302938	1.773252	2.652732
2000-01-02	-2.012384	-0.152303	-0.252590	0.971645
2000-01-03	-1.193222	0.568376	-0.096553	0.339151
2000-01-04	1.438393	0.040858	-1.978462	0.488624
2000-01-05	0.492501	1.287490	0.312481	-0.268951

Item2

	A	B	C	D
2000-01-01	-1.192723	-0.825759	-1.474287	-0.146466
2000-01-02	-0.597677	-1.108512	0.673023	-0.466862
2000-01-03	0.086895	0.120188	1.399064	-1.806952
2000-01-04	0.859977	0.018206	1.295162	-0.807180
2000-01-05	0.963918	-1.076154	0.346418	0.447741

### 6.7.2 iterrows

New in v0.7 is the ability to iterate efficiently through rows of a DataFrame. It returns an iterator yielding each index value along with a Series containing the data in each row:

```
In [139]: for row_index, row in df2.iterrows():
.....:     print '%s\n%s' % (row_index, row)
.....:
```

a

one	0.607032
two	-0.935367

Name: a

b

one	-1.149464
two	0.262314

Name: b

c

one	0.542432
-----	----------



```
two      0.673052
Name: c
```

For instance, a contrived way to transpose the dataframe would be:

```
In [140]: df2 = DataFrame({'x': [1, 2, 3], 'y': [4, 5, 6]})
```

```
In [141]: print df2
```

```
   x  y
0  1  4
1  2  5
2  3  6
```

```
In [142]: print df2.T
```

```
   0  1  2
x  1  2  3
y  4  5  6
```

```
In [143]: df2_t = DataFrame(dict((idx, values) for idx, values in df2.iterrows()))
```

```
In [144]: print df2_t
```

```
   0  1  2
x  1  2  3
y  4  5  6
```

### 6.7.3 itertuples

This method will return an iterator yielding a tuple for each row in the DataFrame. The first element of the tuple will be the row's corresponding index value, while the remaining values are the row values proper.

For instance,

```
In [145]: for r in df2.itertuples(): print r
```

```
(0, 1, 4)
(1, 2, 5)
(2, 3, 6)
```

## 6.8 Vectorized string methods

Series is equipped (as of pandas 0.8.1) with a set of string processing methods that make it easy to operate on each element of the array. Perhaps most importantly, these methods exclude missing/NA values automatically. These are accessed via the Series's `str` attribute and generally have names matching the equivalent (scalar) build-in string methods:

```
In [146]: s = Series(['A', 'B', 'C', 'Aaba', 'Baca', np.nan, 'CABA', 'dog', 'cat'])
```

```
In [147]: s.str.lower()
```

```
Out[147]:
0      a
1      b
2      c
3    aaba
4    baca
5     NaN
6    caba
```

```
7      dog
8      cat

In [148]: s.str.upper()
Out[148]:
0      A
1      B
2      C
3  AABA
4  BACA
5   NaN
6  CABA
7   DOG
8   CAT
```

```
In [149]: s.str.len()
Out[149]:
0      1
1      1
2      1
3      4
4      4
5   NaN
6      4
7      3
8      3
```

Methods like `split` return a Series of lists:

```
In [150]: s2 = Series(['a_b_c', 'c_d_e', np.nan, 'f_g_h'])

In [151]: s2.str.split('_')
Out[151]:
0      ['a', 'b', 'c']
1      ['c', 'd', 'e']
2      NaN
3      ['f', 'g', 'h']
```

Elements in the split lists can be accessed using `get` or `[]` notation:

```
In [152]: s2.str.split('_').str.get(1)
Out[152]:
0      b
1      d
2   NaN
3      g
```

```
In [153]: s2.str.split('_').str[1]
Out[153]:
0      b
1      d
2   NaN
3      g
```

Methods like `replace` and `findall` take regular expressions, too:

```
In [154]: s3 = Series(['A', 'B', 'C', 'Aaba', 'Baca',
.....:                '', np.nan, 'CABA', 'dog', 'cat'])
.....:
```

```
In [155]: s3
```

```
Out[155]:
```

```
0      A
1      B
2      C
3  Aaba
4  Baca
5
6     NaN
7  CABA
8   dog
9   cat
```

```
In [156]: s3.str.replace('^a|dog', 'XX-XX ', case=False)
```

```
Out[156]:
```

```
0      A
1      B
2      C
3  XX-XX ba
4  XX-XX ca
5
6     NaN
7  XX-XX BA
8   XX-XX
9  XX-XX t
```

Method	Description
cat	Concatenate strings
split	Split strings on delimiter
get	Index into each element (retrieve i-th element)
join	Join strings in each element of the Series with passed separator
contains	Return boolean array if each string contains pattern/regex
replace	Replace occurrences of pattern/regex with some other string
repeat	Duplicate values ( <code>s.str.repeat(3)</code> equivalent to <code>x * 3</code> )
pad	Add whitespace to left, right, or both sides of strings
center	Equivalent to <code>pad(side='both')</code>
slice	Slice each string in the Series
slice_replace	Replace slice in each string with passed value
count	Count occurrences of pattern
startswith	Equivalent to <code>str.startswith(pat)</code> for each element
endswith	Equivalent to <code>str.endswith(pat)</code> for each element
findall	Compute list of all occurrences of pattern/regex for each string
match	Call <code>re.match</code> on each element, returning matched groups as list
len	Compute string lengths
strip	Equivalent to <code>str.strip</code>
rstrip	Equivalent to <code>str.rstrip</code>
lstrip	Equivalent to <code>str.lstrip</code>
lower	Equivalent to <code>str.lower</code>
upper	Equivalent to <code>str.upper</code>

## 6.9 Sorting by index and value

There are two obvious kinds of sorting that you may be interested in: sorting by label and sorting by actual values. The primary method for sorting axis labels (indexes) across data structures is the `sort_index` method.

```
In [157]: unsorted_df = df.reindex(index=['a', 'd', 'c', 'b'],
.....:                               columns=['three', 'two', 'one'])
.....:
```

```
In [158]: unsorted_df.sort_index()
Out[158]:
```

	three	two	one
a	NaN	-0.478462	0.588637
b	0.085539	0.719219	-1.167858
c	0.910929	1.129957	0.524038
d	-1.274663	0.739634	NaN

```
In [159]: unsorted_df.sort_index(ascending=False)
Out[159]:
```

	three	two	one
d	-1.274663	0.739634	NaN
c	0.910929	1.129957	0.524038
b	0.085539	0.719219	-1.167858
a	NaN	-0.478462	0.588637

```
In [160]: unsorted_df.sort_index(axis=1)
Out[160]:
```

	one	three	two
a	0.588637	NaN	-0.478462
d	NaN	-1.274663	0.739634
c	0.524038	0.910929	1.129957
b	-1.167858	0.085539	0.719219

`DataFrame.sort_index` can accept an optional `by` argument for `axis=0` which will use an arbitrary vector or a column name of the `DataFrame` to determine the sort order:

```
In [161]: df.sort_index(by='two')
Out[161]:
```

	one	three	two
a	0.588637	NaN	-0.478462
b	-1.167858	0.085539	0.719219
d	NaN	-1.274663	0.739634
c	0.524038	0.910929	1.129957

The `by` argument can take a list of column names, e.g.:

```
In [162]: df = DataFrame({'one': [2, 1, 1, 1], 'two': [1, 3, 2, 4], 'three': [5, 4, 3, 2]})
```

```
In [163]: df[['one', 'two', 'three']].sort_index(by=['one', 'two'])
Out[163]:
```

	one	two	three
2	1	2	3
1	1	3	4
3	1	4	2
0	2	1	5

`Series` has the method `order` (analogous to R's `order` function) which sorts by value, with special treatment of NA values via the `na_last` argument:

```
In [164]: s[2] = np.nan
```

```
In [165]: s.order()
```

```
Out[165]:
```

```
0      A
3    Aaba
1      B
4    Baca
6    CABA
8     cat
7     dog
2     NaN
5     NaN
```

```
In [166]: s.order(na_last=False)
```

```
Out[166]:
```

```
2     NaN
5     NaN
0      A
3    Aaba
1      B
4    Baca
6    CABA
8     cat
7     dog
```

Some other sorting notes / nuances:

- `Series.sort` sorts a `Series` by value in-place. This is to provide compatibility with NumPy methods which expect the `ndarray.sort` behavior.
- `DataFrame.sort` takes a `column` argument instead of `by`. This method will likely be deprecated in a future release in favor of just using `sort_index`.

## 6.10 Copying, type casting

The `copy` method on pandas objects copies the underlying data (though not the axis indexes, since they are immutable) and returns a new object. Note that **it is seldom necessary to copy objects**. For example, there are only a handful of ways to alter a `DataFrame` *in-place*:

- Inserting, deleting, or modifying a column
- Assigning to the `index` or `columns` attributes
- For homogeneous data, directly modifying the values via the `values` attribute or advanced indexing

To be clear, no pandas methods have the side effect of modifying your data; almost all methods return new objects, leaving the original object untouched. If data is modified, it is because you did so explicitly.

Data can be explicitly cast to a NumPy dtype by using the `astype` method or alternately passing the `dtype` keyword argument to the object constructor.

```
In [167]: df = DataFrame(np.arange(12).reshape((4, 3)))
```

```
In [168]: df[0].dtype
```

```
Out[168]: dtype('int64')
```

```
In [169]: df.astype(float)[0].dtype
```

```
Out[169]: dtype('float64')
```

```
In [170]: df = DataFrame(np.arange(12).reshape((4, 3)), dtype=float)
```

```
In [171]: df[0].dtype
```

```
Out[171]: dtype('float64')
```

### 6.10.1 Inferring better types for object columns

The `convert_objects` DataFrame method will attempt to convert `dtype=object` columns to a better NumPy dtype. Occasionally (after transposing multiple times, for example), a mixed-type DataFrame will end up with everything as `dtype=object`. This method attempts to fix that:

```
In [172]: df = DataFrame(randn(6, 3), columns=['a', 'b', 'c'])
```

```
In [173]: df['d'] = 'foo'
```

```
In [174]: df
```

```
Out[174]:
```

	a	b	c	d
0	-1.587663	-0.678060	0.725873	foo
1	-1.214581	0.170874	-0.859984	foo
2	-0.736275	-0.596860	1.321141	foo
3	-0.042158	0.179324	0.858098	foo
4	-1.082848	0.040125	0.700082	foo
5	-1.331745	0.729971	-0.383658	foo

```
In [175]: df = df.T.T
```

```
In [176]: df.dtypes
```

```
Out[176]:
```

a	object
b	object
c	object
d	object

```
In [177]: converted = df.convert_objects()
```

```
In [178]: converted.dtypes
```

```
Out[178]:
```

a	float64
b	float64
c	float64
d	object

## 6.11 Pickling and serialization

All pandas objects are equipped with `save` methods which use Python's `cPickle` module to save data structures to disk using the pickle format.

```
In [179]: df
```

```
Out[179]:
```

	a	b	c	d
0	-1.587663	-0.6780604	0.725873	foo

```

1  -1.214581    0.1708745 -0.8599843  foo
2  -0.7362747 -0.5968605  1.321141   foo
3  -0.04215821  0.1793244  0.8580981  foo
4  -1.082848   0.04012511  0.7000821  foo
5  -1.331745   0.7299706 -0.3836576  foo

```

```
In [180]: df.save('foo.pickle')
```

The load function in the pandas namespace can be used to load any pickled pandas object (or any other pickled object) from file:

```
In [181]: load('foo.pickle')
```

```
Out[181]:
```

	a	b	c	d
0	-1.587663	-0.6780604	0.725873	foo
1	-1.214581	0.1708745	-0.8599843	foo
2	-0.7362747	-0.5968605	1.321141	foo
3	-0.04215821	0.1793244	0.8580981	foo
4	-1.082848	0.04012511	0.7000821	foo
5	-1.331745	0.7299706	-0.3836576	foo

There is also a save function which takes any object as its first argument:

```
In [182]: save(df, 'foo.pickle')
```

```
In [183]: load('foo.pickle')
```

```
Out[183]:
```

	a	b	c	d
0	-1.587663	-0.6780604	0.725873	foo
1	-1.214581	0.1708745	-0.8599843	foo
2	-0.7362747	-0.5968605	1.321141	foo
3	-0.04215821	0.1793244	0.8580981	foo
4	-1.082848	0.04012511	0.7000821	foo
5	-1.331745	0.7299706	-0.3836576	foo

## 6.12 Console Output Formatting

Use the `set_eng_float_format` function in the `pandas.core.common` module to alter the floating-point formatting of pandas objects to produce a particular format.

For instance:

```
In [184]: set_eng_float_format(accuracy=3, use_eng_prefix=True)
```

```
In [185]: df['a']/1.e3
```

```
Out[185]:
```

0	-1.588m
1	-1.215m
2	-736.275u
3	-42.158u
4	-1.083m
5	-1.332m

Name: a

```
In [186]: df['a']/1.e6
```

```
Out[186]:
```

0	-1.588u
---	---------

```
1      -1.215u
2     -736.275n
3     -42.158n
4      -1.083u
5      -1.332u
Name: a
```

The `set_printoptions` function has a number of options for controlling how floating point numbers are formatted (using the `precision` argument) in the console and `.`. The `max_rows` and `max_columns` control how many rows and columns of `DataFrame` objects are shown by default. If `max_columns` is set to 0 (the default, in fact), the library will attempt to fit the `DataFrame`'s string representation into the current terminal width, and defaulting to the summary view otherwise.



# INDEXING AND SELECTING DATA

The axis labeling information in pandas objects serves many purposes:

- Identifies data (i.e. provides *metadata*) using known indicators, important for analysis, visualization, and interactive console display
- Enables automatic and explicit data alignment
- Allows intuitive getting and setting of subsets of the data set

In this section / chapter, we will focus on the final point: namely, how to slice, dice, and generally get and set subsets of pandas objects. The primary focus will be on Series and DataFrame as they have received more development attention in this area. Expect more work to be invested higher-dimensional data structures (including Panel) in the future, especially in label-based advanced indexing.

## 7.1 Basics

As mentioned when introducing the data structures in the [last section](#), the primary function of indexing with `[]` (a.k.a. `__getitem__` for those familiar with implementing class behavior in Python) is selecting out lower-dimensional slices. Thus,

- **Series:** `series[label]` returns a scalar value
- **DataFrame:** `frame[colname]` returns a Series corresponding to the passed column name
- **Panel:** `panel[itemname]` returns a DataFrame corresponding to the passed item name

Here we construct a simple time series data set to use for illustrating the indexing functionality:

```
In [542]: dates = np.asarray(date_range('1/1/2000', periods=8))
```

```
In [543]: df = DataFrame(randn(8, 4), index=dates, columns=['A', 'B', 'C', 'D'])
```

```
In [544]: df
```

```
Out[544]:
```

	A	B	C	D
2000-01-01	0.469112	-0.282863	-1.509059	-1.135632
2000-01-02	1.212112	-0.173215	0.119209	-1.044236
2000-01-03	-0.861849	-2.104569	-0.494929	1.071804
2000-01-04	0.721555	-0.706771	-1.039575	0.271860
2000-01-05	-0.424972	0.567020	0.276232	-1.087401
2000-01-06	-0.673690	0.113648	-1.478427	0.524988
2000-01-07	0.404705	0.577046	-1.715002	-1.039268
2000-01-08	-0.370647	-1.157892	-1.344312	0.844885

```
In [545]: panel = Panel({'one' : df, 'two' : df - df.mean()})
```

```
In [546]: panel
```

```
Out[546]:  
<class 'pandas.core.panel.Panel'>  
Dimensions: 2 (items) x 8 (major) x 4 (minor)  
Items: one to two  
Major axis: 2000-01-01 00:00:00 to 2000-01-08 00:00:00  
Minor axis: A to D
```

---

**Note:** None of the indexing functionality is time series specific unless specifically stated.

---

Thus, as per above, we have the most basic indexing using []:

```
In [547]: s = df['A']
```

```
In [548]: s[dates[5]]
```

```
Out[548]: -0.67368970808837059
```

```
In [549]: panel['two']
```

```
Out[549]:
```

	A	B	C	D
2000-01-01	0.409571	0.113086	-0.610826	-0.936507
2000-01-02	1.152571	0.222735	1.017442	-0.845111
2000-01-03	-0.921390	-1.708620	0.403304	1.270929
2000-01-04	0.662014	-0.310822	-0.141342	0.470985
2000-01-05	-0.484513	0.962970	1.174465	-0.888276
2000-01-06	-0.733231	0.509598	-0.580194	0.724113
2000-01-07	0.345164	0.972995	-0.816769	-0.840143
2000-01-08	-0.430188	-0.761943	-0.446079	1.044010

## 7.1.1 Fast scalar value getting and setting

Since indexing with [] must handle a lot of cases (single-label access, slicing, boolean indexing, etc.), it has a bit of overhead in order to figure out what you're asking for. If you only want to access a scalar value, the fastest way is to use the `get_value` method, which is implemented on all of the data structures:

```
In [550]: s.get_value(dates[5])
```

```
Out[550]: -0.67368970808837059
```

```
In [551]: df.get_value(dates[5], 'A')
```

```
Out[551]: -0.67368970808837059
```

There is an analogous `set_value` method which has the additional capability of enlarging an object. This method *always* returns a reference to the object it modified, which in the fast of enlargement, will be a **new object**:

```
In [552]: df.set_value(dates[5], 'E', 7)
```

```
Out[552]:
```

	A	B	C	D	E
2000-01-01	0.469112	-0.282863	-1.509059	-1.135632	NaN
2000-01-02	1.212112	-0.173215	0.119209	-1.044236	NaN
2000-01-03	-0.861849	-2.104569	-0.494929	1.071804	NaN
2000-01-04	0.721555	-0.706771	-1.039575	0.271860	NaN
2000-01-05	-0.424972	0.567020	0.276232	-1.087401	NaN
2000-01-06	-0.673690	0.113648	-1.478427	0.524988	7

```
2000-01-07    0.404705    0.577046 -1.715002 -1.039268 NaN
2000-01-08   -0.370647   -1.157892 -1.344312    0.844885 NaN
```

## 7.1.2 Additional Column Access

You may access a column on a dataframe directly as an attribute:

```
In [553]: df.A
Out[553]:
2000-01-01    0.469112
2000-01-02    1.212112
2000-01-03   -0.861849
2000-01-04    0.721555
2000-01-05   -0.424972
2000-01-06   -0.673690
2000-01-07    0.404705
2000-01-08   -0.370647
Name: A
```

If you are using the IPython environment, you may also use tab-completion to see the accessible columns of a DataFrame.

You can pass a list of columns to `[]` to select columns in that order: If a column is not contained in the DataFrame, an exception will be raised. Multiple columns can also be set in this manner:

```
In [554]: df
Out[554]:
```

	A	B	C	D
2000-01-01	0.469112	-0.282863	-1.509059	-1.135632
2000-01-02	1.212112	-0.173215	0.119209	-1.044236
2000-01-03	-0.861849	-2.104569	-0.494929	1.071804
2000-01-04	0.721555	-0.706771	-1.039575	0.271860
2000-01-05	-0.424972	0.567020	0.276232	-1.087401
2000-01-06	-0.673690	0.113648	-1.478427	0.524988
2000-01-07	0.404705	0.577046	-1.715002	-1.039268
2000-01-08	-0.370647	-1.157892	-1.344312	0.844885

```
In [555]: df[['B', 'A']] = df[['A', 'B']]
```

```
In [556]: df
Out[556]:
```

	A	B	C	D
2000-01-01	-0.282863	0.469112	-1.509059	-1.135632
2000-01-02	-0.173215	1.212112	0.119209	-1.044236
2000-01-03	-2.104569	-0.861849	-0.494929	1.071804
2000-01-04	-0.706771	0.721555	-1.039575	0.271860
2000-01-05	0.567020	-0.424972	0.276232	-1.087401
2000-01-06	0.113648	-0.673690	-1.478427	0.524988
2000-01-07	0.577046	0.404705	-1.715002	-1.039268
2000-01-08	-1.157892	-0.370647	-1.344312	0.844885

You may find this useful for applying a transform (in-place) to a subset of the columns.

## 7.1.3 Data slices on other axes

It's certainly possible to retrieve data slices along the other axes of a DataFrame or Panel. We tend to refer to these slices as *cross-sections*. DataFrame has the `xs` function for retrieving rows as Series and Panel has the analogous

`major_xs` and `minor_xs` functions for retrieving slices as DataFrames for a given `major_axis` or `minor_axis` label, respectively.

```
In [557]: date = dates[5]
```

```
In [558]: df.xs(date)
```

```
Out[558]:
A    0.113648
B   -0.673690
C   -1.478427
D    0.524988
Name: 2000-01-06 00:00:00
```

```
In [559]: panel.major_xs(date)
```

```
Out[559]:
           one      two
A -0.673690 -0.733231
B  0.113648  0.509598
C -1.478427 -0.580194
D  0.524988  0.724113
```

```
In [560]: panel.minor_xs('A')
```

```
Out[560]:
           one      two
2000-01-01  0.469112  0.409571
2000-01-02  1.212112  1.152571
2000-01-03 -0.861849 -0.921390
2000-01-04  0.721555  0.662014
2000-01-05 -0.424972 -0.484513
2000-01-06 -0.673690 -0.733231
2000-01-07  0.404705  0.345164
2000-01-08 -0.370647 -0.430188
```

### 7.1.4 Slicing ranges

The most robust and consistent way of slicing ranges along arbitrary axes is described in the [Advanced indexing](#) section detailing the `.ix` method. For now, we explain the semantics of slicing using the `[]` operator.

With Series, the syntax works exactly as with an ndarray, returning a slice of the values and the corresponding labels:

```
In [561]: s[:5]
```

```
Out[561]:
2000-01-01    -0.282863
2000-01-02    -0.173215
2000-01-03    -2.104569
2000-01-04    -0.706771
2000-01-05     0.567020
Name: A
```

```
In [562]: s[::2]
```

```
Out[562]:
2000-01-01    -0.282863
2000-01-03    -2.104569
2000-01-05     0.567020
2000-01-07     0.577046
Name: A
```

```
In [563]: s[::-1]
```

```
Out [563]:
2000-01-08    -1.157892
2000-01-07     0.577046
2000-01-06     0.113648
2000-01-05     0.567020
2000-01-04    -0.706771
2000-01-03    -2.104569
2000-01-02    -0.173215
2000-01-01    -0.282863
Name: A
```

Note that setting works as well:

```
In [564]: s2 = s.copy()
```

```
In [565]: s2[:5] = 0
```

```
In [566]: s2
Out [566]:
2000-01-01     0.000000
2000-01-02     0.000000
2000-01-03     0.000000
2000-01-04     0.000000
2000-01-05     0.000000
2000-01-06     0.113648
2000-01-07     0.577046
2000-01-08    -1.157892
Name: A
```

With DataFrame, slicing inside of `[]` **slices the rows**. This is provided largely as a convenience since it is such a common operation.

```
In [567]: df[:3]
Out [567]:
```

	A	B	C	D
2000-01-01	-0.282863	0.469112	-1.509059	-1.135632
2000-01-02	-0.173215	1.212112	0.119209	-1.044236
2000-01-03	-2.104569	-0.861849	-0.494929	1.071804

```
In [568]: df[::-1]
Out [568]:
```

	A	B	C	D
2000-01-08	-1.157892	-0.370647	-1.344312	0.844885
2000-01-07	0.577046	0.404705	-1.715002	-1.039268
2000-01-06	0.113648	-0.673690	-1.478427	0.524988
2000-01-05	0.567020	-0.424972	0.276232	-1.087401
2000-01-04	-0.706771	0.721555	-1.039575	0.271860
2000-01-03	-2.104569	-0.861849	-0.494929	1.071804
2000-01-02	-0.173215	1.212112	0.119209	-1.044236
2000-01-01	-0.282863	0.469112	-1.509059	-1.135632

## 7.1.5 Boolean indexing

Another common operation is the use of boolean vectors to filter the data.

Using a boolean vector to index a Series works exactly as in a numpy ndarray:

```
In [569]: s[s > 0]
Out[569]:
2000-01-05    0.567020
2000-01-06    0.113648
2000-01-07    0.577046
Name: A
```

```
In [570]: s[(s < 0) & (s > -0.5)]
Out[570]:
2000-01-01   -0.282863
2000-01-02   -0.173215
Name: A
```

You may select rows from a DataFrame using a boolean vector the same length as the DataFrame's index (for example, something derived from one of the columns of the DataFrame):

```
In [571]: df[df['A'] > 0]
Out[571]:
```

	A	B	C	D
2000-01-05	0.567020	-0.424972	0.276232	-1.087401
2000-01-06	0.113648	-0.673690	-1.478427	0.524988
2000-01-07	0.577046	0.404705	-1.715002	-1.039268

Consider the `isin` method of Series, which returns a boolean vector that is true wherever the Series elements exist in the passed list. This allows you to select rows where one or more columns have values you want:

```
In [572]: df2 = DataFrame({'a' : ['one', 'one', 'two', 'three', 'two', 'one', 'six'],
.....:                  'b' : ['x', 'y', 'y', 'x', 'y', 'x', 'x'],
.....:                  'c' : randn(7)})
.....:
```

```
In [573]: df2[df2['a'].isin(['one', 'two'])]
Out[573]:
```

	a	b	c
0	one	x	1.075770
1	one	y	-0.109050
2	two	y	1.643563
4	two	y	0.357021
5	one	x	-0.674600

List comprehensions and map method of Series can also be used to produce more complex criteria:

```
# only want 'two' or 'three'
In [574]: criterion = df2['a'].map(lambda x: x.startswith('t'))
```

```
In [575]: df2[criterion]
Out[575]:
```

	a	b	c
2	two	y	1.643563
3	three	x	-1.469388
4	two	y	0.357021

```
# equivalent but slower
```

```
In [576]: df2[[x.startswith('t') for x in df2['a']]]
Out[576]:
```

	a	b	c
2	two	y	1.643563
3	three	x	-1.469388
4	two	y	0.357021

```
# Multiple criteria
In [577]: df2[criterion & (df2['b'] == 'x')]
Out[577]:
      a  b      c
3  three  x -1.469388
```

Note, with the *advanced indexing* `ix` method, you may select along more than one axis using boolean vectors combined with other indexing expressions.

## 7.1.6 Indexing a DataFrame with a boolean DataFrame

You may wish to set values on a DataFrame based on some boolean criteria derived from itself or another DataFrame or set of DataFrames. This can be done intuitively like so:

```
In [578]: df2 = df.copy()

In [579]: df2 < 0
Out[579]:
      A      B      C      D
2000-01-01  True  False  True  True
2000-01-02  True  False  False  True
2000-01-03  True   True   True  False
2000-01-04  True  False  True  False
2000-01-05  False  True  False  True
2000-01-06  False  True   True  False
2000-01-07  False  False  True  True
2000-01-08  True   True   True  False
```

```
In [580]: df2[df2 < 0] = 0
```

```
In [581]: df2
Out[581]:
      A      B      C      D
2000-01-01  0.000000  0.469112  0.000000  0.000000
2000-01-02  0.000000  1.212112  0.119209  0.000000
2000-01-03  0.000000  0.000000  0.000000  1.071804
2000-01-04  0.000000  0.721555  0.000000  0.271860
2000-01-05  0.567020  0.000000  0.276232  0.000000
2000-01-06  0.113648  0.000000  0.000000  0.524988
2000-01-07  0.577046  0.404705  0.000000  0.000000
2000-01-08  0.000000  0.000000  0.000000  0.844885
```

Note that such an operation requires that the boolean DataFrame is indexed exactly the same.

## 7.1.7 Take Methods

Similar to numpy ndarrays, pandas Index, Series, and DataFrame also provides the `take` method that retrieves elements along a given axis at the given indices. The given indices must be either a list or an ndarray of integer index positions.

```
In [582]: index = Index(randint(0, 1000, 10))

In [583]: index
Out[583]: Int64Index([969, 412, 496, 195, 288, 101, 881, 900, 732, 658])

In [584]: positions = [0, 9, 3]
```

```
In [585]: index[positions]
Out[585]: Int64Index([969, 658, 195])
```

```
In [586]: index.take(positions)
Out[586]: Int64Index([969, 658, 195])
```

```
In [587]: ser = Series(randn(10))
```

```
In [588]: ser.ix[positions]
Out[588]:
0    -0.968914
9    -1.131345
3     1.247642
```

```
In [589]: ser.take(positions)
Out[589]:
0    -0.968914
9    -1.131345
3     1.247642
```

For DataFrames, the given indices should be a 1d list or ndarray that specifies row or column positions.

```
In [590]: frm = DataFrame(randn(5, 3))
```

```
In [591]: frm.take([1, 4, 3])
Out[591]:
```

	0	1	2
1	-0.932132	1.956030	0.017587
4	-0.077118	-0.408530	-0.862495
3	-1.143704	0.215897	1.193555

```
In [592]: frm.take([0, 2], axis=1)
Out[592]:
```

	0	2
0	-0.089329	-0.945867
1	-0.932132	0.017587
2	-0.016692	0.254161
3	-1.143704	1.193555
4	-0.077118	-0.862495

It is important to note that the `take` method on pandas objects are not intended to work on boolean indices and may return unexpected results.

```
In [593]: arr = randn(10)
```

```
In [594]: arr.take([False, False, True, True])
Out[594]: array([ 1.3461,  1.3461,  1.5118,  1.5118])
```

```
In [595]: arr[[0, 1]]
Out[595]: array([ 1.3461,  1.5118])
```

```
In [596]: ser = Series(randn(10))
```

```
In [597]: ser.take([False, False, True, True])
Out[597]:
0    -0.105381
0    -0.105381
1    -0.532532
```



```
1    -0.532532
```

```
In [598]: ser.ix[[0, 1]]
```

```
Out[598]:
```

```
0    -0.105381
```

```
1    -0.532532
```

Finally, as a small note on performance, because the `take` method handles a narrower range of inputs, it can offer performance that is a good deal faster than fancy indexing.

## 7.1.8 Duplicate Data

If you want to identify and remove duplicate rows in a `DataFrame`, there are two methods that will help: `duplicated` and `drop_duplicates`. Each takes as an argument the columns to use to identify duplicated rows.

`duplicated` returns a boolean vector whose length is the number of rows, and which indicates whether a row is duplicated.

`drop_duplicates` removes duplicate rows.

By default, the first observed row of a duplicate set is considered unique, but each method has a `take_last` parameter that indicates the last observed row should be taken instead.

```
In [599]: df2 = DataFrame({'a' : ['one', 'one', 'two', 'three', 'two', 'one', 'six'],
.....:                  'b' : ['x', 'y', 'y', 'x', 'y', 'x', 'x'],
.....:                  'c' : np.random.randn(7)})
.....:
```

```
In [600]: df2.duplicated(['a', 'b'])
```

```
Out[600]:
```

```
0    False
```

```
1    False
```

```
2    False
```

```
3    False
```

```
4     True
```

```
5     True
```

```
6    False
```

```
In [601]: df2.drop_duplicates(['a', 'b'])
```

```
Out[601]:
```

```
   a  b      c
0  one x -0.339355
1  one y  0.593616
2  two y  0.884345
3 three x  1.591431
6  six x  0.435589
```

```
In [602]: df2.drop_duplicates(['a', 'b'], take_last=True)
```

```
Out[602]:
```

```
   a  b      c
1  one y  0.593616
3 three x  1.591431
4  two y  0.141809
5  one x  0.220390
6  six x  0.435589
```

### 7.1.9 Dictionary-like get method

Each of Series, DataFrame, and Panel have a `get` method which can return a default value.

```
In [603]: s = Series([1,2,3], index=['a','b','c'])

In [604]: s.get('a')                # equivalent to s['a']
Out[604]: 1

In [605]: s.get('x', default=-1)
Out[605]: -1
```

## 7.2 Advanced indexing with labels

We have avoided excessively overloading the `[]` / `__getitem__` operator to keep the basic functionality of the pandas objects straightforward and simple. However, there are often times when you may wish get a subset (or analogously set a subset) of the data in a way that is not straightforward using the combination of `reindex` and `[]`. Complicated setting operations are actually quite difficult because `reindex` usually returns a copy.

By *advanced* indexing we are referring to a special `.ix` attribute on pandas objects which enable you to do getting/setting operations on a DataFrame, for example, with matrix/ndarray-like semantics. Thus you can combine the following kinds of indexing:

- An integer or single label, e.g. 5 or 'a'
- A list or array of labels ['a', 'b', 'c'] or integers [4, 3, 0]
- A slice object with ints 1:7 or labels 'a':'f'
- A boolean array

We'll illustrate all of these methods. First, note that this provides a concise way of reindexing on multiple axes at once:

```
In [606]: subindex = dates[[3,4,5]]

In [607]: df.reindex(index=subindex, columns=['C', 'B'])
Out[607]:
```

	C	B
2000-01-04	-1.039575	0.721555
2000-01-05	0.276232	-0.424972
2000-01-06	-1.478427	-0.673690

```
In [608]: df.ix[subindex, ['C', 'B']]
Out[608]:
```

	C	B
2000-01-04	-1.039575	0.721555
2000-01-05	0.276232	-0.424972
2000-01-06	-1.478427	-0.673690

Assignment / setting values is possible when using `ix`:

```
In [609]: df2 = df.copy()

In [610]: df2.ix[subindex, ['C', 'B']] = 0

In [611]: df2
Out[611]:
```

	A	B	C	D
2000-01-01	-0.282863	0.469112	-1.509059	-1.135632

```

2000-01-02 -0.173215  1.212112  0.119209 -1.044236
2000-01-03 -2.104569 -0.861849 -0.494929  1.071804
2000-01-04 -0.706771  0.000000  0.000000  0.271860
2000-01-05  0.567020  0.000000  0.000000 -1.087401
2000-01-06  0.113648  0.000000  0.000000  0.524988
2000-01-07  0.577046  0.404705 -1.715002 -1.039268
2000-01-08 -1.157892 -0.370647 -1.344312  0.844885

```

Indexing with an array of integers can also be done:

```
In [612]: df.ix[[4,3,1]]
```

```
Out [612]:
```

	A	B	C	D
2000-01-05	0.567020	-0.424972	0.276232	-1.087401
2000-01-04	-0.706771	0.721555	-1.039575	0.271860
2000-01-02	-0.173215	1.212112	0.119209	-1.044236

```
In [613]: df.ix[dates[[4,3,1]]]
```

```
Out [613]:
```

	A	B	C	D
2000-01-05	0.567020	-0.424972	0.276232	-1.087401
2000-01-04	-0.706771	0.721555	-1.039575	0.271860
2000-01-02	-0.173215	1.212112	0.119209	-1.044236

Slicing has standard Python semantics for integer slices:

```
In [614]: df.ix[1:7, :2]
```

```
Out [614]:
```

	A	B
2000-01-02	-0.173215	1.212112
2000-01-03	-2.104569	-0.861849
2000-01-04	-0.706771	0.721555
2000-01-05	0.567020	-0.424972
2000-01-06	0.113648	-0.673690
2000-01-07	0.577046	0.404705

Slicing with labels is semantically slightly different because the slice start and stop are **inclusive** in the label-based case:

```
In [615]: date1, date2 = dates[[2, 4]]
```

```
In [616]: print date1, date2
```

```
1970-01-11 232:00:00 1970-01-11 24:00:00
```

```
In [617]: df.ix[date1:date2]
```

```
Out [617]:
```

Empty DataFrame  
Columns: array([A, B, C, D], dtype=object)  
Index: <class 'pandas.tseries.index.DatetimeIndex'>  
Length: 0, Freq: None, Timezone: None

```
In [618]: df['A'].ix[date1:date2]
```

```
Out [618]: TimeSeries([], dtype=float64)
```

Getting and setting rows in a DataFrame, especially by their location, is much easier:

```
In [619]: df2 = df[:5].copy()
```

```
In [620]: df2.ix[3]
```

```
Out [620]:
```

```
A    -0.706771
B     0.721555
C    -1.039575
D     0.271860
Name: 2000-01-04 00:00:00
```

```
In [621]: df2.ix[3] = np.arange(len(df2.columns))
```

```
In [622]: df2
```

```
Out [622]:
```

	A	B	C	D
2000-01-01	-0.282863	0.469112	-1.509059	-1.135632
2000-01-02	-0.173215	1.212112	0.119209	-1.044236
2000-01-03	-2.104569	-0.861849	-0.494929	1.071804
2000-01-04	0.000000	1.000000	2.000000	3.000000
2000-01-05	0.567020	-0.424972	0.276232	-1.087401

Column or row selection can be combined as you would expect with arrays of labels or even boolean vectors:

```
In [623]: df.ix[df['A'] > 0, 'B']
```

```
Out [623]:
```

2000-01-05	-0.424972
2000-01-06	-0.673690
2000-01-07	0.404705

Name: B

```
In [624]: df.ix[date1:date2, 'B']
```

```
Out [624]: TimeSeries([], dtype=float64)
```

```
In [625]: df.ix[date1, 'B']
```

```
Out [625]: -0.86184896334779992
```

Slicing with labels is closely related to the `truncate` method which does precisely `.ix[start:stop]` but returns a copy (for legacy reasons).

## 7.2.1 Returning a view versus a copy

The rules about when a view on the data is returned are entirely dependent on NumPy. Whenever an array of labels or a boolean vector are involved in the indexing operation, the result will be a copy. With single label / scalar indexing and slicing, e.g. `df.ix[3:6]` or `df.ix[:, 'A']`, a view will be returned.

## 7.2.2 The `select` method

Another way to extract slices from an object is with the `select` method of Series, DataFrame, and Panel. This method should be used only when there is no more direct way. `select` takes a function which operates on labels along `axis` and returns a boolean. For instance:

```
In [626]: df.select(lambda x: x == 'A', axis=1)
```

```
Out [626]:
```

	A
2000-01-01	-0.282863
2000-01-02	-0.173215
2000-01-03	-2.104569
2000-01-04	-0.706771
2000-01-05	0.567020
2000-01-06	0.113648

```
2000-01-07    0.577046
2000-01-08   -1.157892
```

### 7.2.3 The lookup method

Sometimes you want to extract a set of values given a sequence of row labels and column labels, and the `lookup` method allows for this and returns a numpy array. For instance,

```
In [627]: dflookup = DataFrame(np.random.rand(20,4), columns = ['A','B','C','D'])

In [628]: dflookup.lookup(xrange(0,10,2), ['B','C','A','B','D'])
Out[628]: array([ 0.0227,  0.4199,  0.529 ,  0.9674,  0.5357])
```

### 7.2.4 Advanced indexing with integer labels

Label-based indexing with integer axis labels is a thorny topic. It has been discussed heavily on mailing lists and among various members of the scientific Python community. In pandas, our general viewpoint is that labels matter more than integer locations. Therefore, with an integer axis index *only* label-based indexing is possible with the standard tools like `.ix`. The following code will generate exceptions:

```
s = Series(range(5))
s[-1]
df = DataFrame(np.random.randn(5, 4))
df
df.ix[-2:]
```

This deliberate decision was made to prevent ambiguities and subtle bugs (many users reported finding bugs when the API change was made to stop “falling back” on position-based indexing).

### 7.2.5 Setting values in mixed-type DataFrame

Setting values on a mixed-type DataFrame or Panel is supported when using scalar values, though setting arbitrary vectors is not yet supported:

```
In [629]: df2 = df[:4]

In [630]: df2['foo'] = 'bar'

In [631]: print df2
```

	A	B	C	D	foo
2000-01-01	-0.282863	0.469112	-1.509059	-1.135632	bar
2000-01-02	-0.173215	1.212112	0.119209	-1.044236	bar
2000-01-03	-2.104569	-0.861849	-0.494929	1.071804	bar
2000-01-04	-0.706771	0.721555	-1.039575	0.271860	bar

```
In [632]: df2.ix[2] = np.nan

In [633]: print df2
```

	A	B	C	D	foo
2000-01-01	-0.282863	0.469112	-1.509059	-1.135632	bar
2000-01-02	-0.173215	1.212112	0.119209	-1.044236	bar
2000-01-03	NaN	NaN	NaN	NaN	NaN
2000-01-04	-0.706771	0.721555	-1.039575	0.271860	bar

```
In [634]: print df2.dtypes
A         float64
B         float64
C         float64
D         float64
foo        object
```

## 7.3 Index objects

The pandas Index class and its subclasses can be viewed as implementing an *ordered set* in addition to providing the support infrastructure necessary for lookups, data alignment, and reindexing. The easiest way to create one directly is to pass a list or other sequence to Index:

```
In [635]: index = Index(['e', 'd', 'a', 'b'])
```

```
In [636]: index
Out[636]: Index([e, d, a, b], dtype=object)
```

```
In [637]: 'd' in index
Out[637]: True
```

You can also pass a name to be stored in the index:

```
In [638]: index = Index(['e', 'd', 'a', 'b'], name='something')
```

```
In [639]: index.name
Out[639]: 'something'
```

Starting with pandas 0.5, the name, if set, will be shown in the console display:

```
In [640]: index = Index(range(5), name='rows')
```

```
In [641]: columns = Index(['A', 'B', 'C'], name='cols')
```

```
In [642]: df = DataFrame(np.random.randn(5, 3), index=index, columns=columns)
```

```
In [643]: df
Out[643]:
```

cols	A	B	C
rows			
0	0.192451	0.629675	-1.425966
1	1.857704	-1.193545	0.677510
2	-0.153931	0.520091	-1.475051
3	0.722570	-0.322646	-1.601631
4	0.778033	-0.289342	0.233141

```
In [644]: df['A']
Out[644]:
```

rows	A
0	0.192451
1	1.857704
2	-0.153931
3	0.722570
4	0.778033

Name: A

### 7.3.1 Set operations on Index objects

The three main operations are `union` (`|`), `intersection` (`&`), and `diff` (`-`). These can be directly called as instance methods or used via overloaded operators:

```
In [645]: a = Index(['c', 'b', 'a'])

In [646]: b = Index(['c', 'e', 'd'])

In [647]: a.union(b)
Out[647]: Index([a, b, c, d, e], dtype=object)

In [648]: a | b
Out[648]: Index([a, b, c, d, e], dtype=object)

In [649]: a & b
Out[649]: Index([c], dtype=object)

In [650]: a - b
Out[650]: Index([a, b], dtype=object)
```

### 7.3.2 `isin` method of Index objects

One additional operation is the `isin` method that works analogously to the `Series.isin` method found [here](#).

## 7.4 Hierarchical indexing (MultiIndex)

Hierarchical indexing (also referred to as “multi-level” indexing) is brand new in the pandas 0.4 release. It is very exciting as it opens the door to some quite sophisticated data analysis and manipulation, especially for working with higher dimensional data. In essence, it enables you to store and manipulate data with an arbitrary number of dimensions in lower dimensional data structures like `Series` (1d) and `DataFrame` (2d).

In this section, we will show what exactly we mean by “hierarchical” indexing and how it integrates with the all of the pandas indexing functionality described above and in prior sections. Later, when discussing *group by* and *pivoting and reshaping data*, we’ll show non-trivial applications to illustrate how it aids in structuring data for analysis.

---

**Note:** Given that hierarchical indexing is so new to the library, it is definitely “bleeding-edge” functionality but is certainly suitable for production. But, there may inevitably be some minor API changes as more use cases are explored and any weaknesses in the design / implementation are identified. pandas aims to be “eminently usable” so any feedback about new functionality like this is extremely helpful.

---

### 7.4.1 Creating a MultiIndex (hierarchical index) object

The `MultiIndex` object is the hierarchical analogue of the standard `Index` object which typically stores the axis labels in pandas objects. You can think of `MultiIndex` an array of tuples where each tuple is unique. A `MultiIndex` can be created from a list of arrays (using `MultiIndex.from_arrays`) or an array of tuples (using `MultiIndex.from_tuples`).

```
In [651]: arrays = [['bar', 'bar', 'baz', 'baz', 'foo', 'foo', 'qux', 'qux'],
.....:             ['one', 'two', 'one', 'two', 'one', 'two', 'one', 'two']]
.....:
```

```
In [652]: tuples = zip(*arrays)
```

```
In [653]: tuples
```

```
Out [653]:  
[('bar', 'one'),  
 ('bar', 'two'),  
 ('baz', 'one'),  
 ('baz', 'two'),  
 ('foo', 'one'),  
 ('foo', 'two'),  
 ('qux', 'one'),  
 ('qux', 'two')]
```

```
In [654]: index = MultiIndex.from_tuples(tuples, names=['first', 'second'])
```

```
In [655]: s = Series(randn(8), index=index)
```

```
In [656]: s
```

```
Out [656]:  
first second  
bar      one    -0.223540  
          two     0.542054  
baz      one    -0.688585  
          two    -0.352676  
foo      one    -0.711411  
          two    -2.122599  
qux      one     1.962935  
          two     1.672027
```

As a convenience, you can pass a list of arrays directly into Series or DataFrame to construct a MultiIndex automatically:

```
In [657]: arrays = [np.array(['bar', 'bar', 'baz', 'baz', 'foo', 'foo', 'qux', 'qux']),  
.....:              np.array(['one', 'two', 'one', 'two', 'one', 'two', 'one', 'two'])]  
.....:
```

```
In [658]: s = Series(randn(8), index=arrays)
```

```
In [659]: s
```

```
Out [659]:  
bar      one    -0.880984  
          two     0.997289  
baz      one    -1.693316  
          two    -0.179129  
foo      one    -1.598062  
          two     0.936914  
qux      one     0.912560  
          two    -1.003401
```

```
In [660]: df = DataFrame(randn(8, 4), index=arrays)
```

```
In [661]: df
```

```
Out [661]:  
          0         1         2         3  
bar one  1.632781 -0.724626  0.178219  0.310610  
      two -0.108002 -0.974226 -1.147708 -2.281374  
baz one   0.760010 -0.742532  1.533318  2.495362  
      two -0.432771 -0.068954  0.043520  0.112246
```



```
foo one    0.871721 -0.816064 -0.784880  1.030659
     two    0.187483 -1.933946  0.377312  0.734122
qux one    2.141616 -0.011225  0.048869 -1.360687
     two   -0.479010 -0.859661 -0.231595 -0.527750
```

All of the `MultiIndex` constructors accept a `names` argument which stores string names for the levels themselves. If no names are provided, some arbitrary ones will be assigned:

```
In [662]: index.names
Out[662]: ['first', 'second']
```

This index can back any axis of a pandas object, and the number of **levels** of the index is up to you:

```
In [663]: df = DataFrame(randn(3, 8), index=['A', 'B', 'C'], columns=index)
```

```
In [664]: df
Out[664]:
<class 'pandas.core.frame.DataFrame'>
Index: 3 entries, A to C
Data columns:
('bar', 'one')      3  non-null values
('bar', 'two')      3  non-null values
('baz', 'one')      3  non-null values
('baz', 'two')      3  non-null values
('foo', 'one')      3  non-null values
('foo', 'two')      3  non-null values
('qux', 'one')      3  non-null values
('qux', 'two')      3  non-null values
dtypes: float64(8)
```

```
In [665]: DataFrame(randn(6, 6), index=index[:6], columns=index[:6])
Out[665]:
first      bar      baz      foo
second     one      two     one      two     one      two
first second
bar   one   -1.993606 -1.927385 -2.027924  1.624972  0.551135  3.059267
      two    0.455264 -0.030740  0.935716  1.061192 -2.107852  0.199905
baz   one    0.323586 -0.641630 -0.587514  0.053897  0.194889 -0.381994
      two    0.318587  2.089075 -0.728293 -0.090255 -0.748199  1.318931
foo   one   -2.029766  0.792652  0.461007 -0.542749 -0.305384 -0.479195
      two    0.095031 -0.270099 -0.707140 -0.773882  0.229453  0.304418
```

We’ve “sparsified” the higher levels of the indexes to make the console output a bit easier on the eyes.

It’s worth keeping in mind that there’s nothing preventing you from using tuples as atomic labels on an axis:

```
In [666]: Series(randn(8), index=tuples)
Out[666]:
('bar', 'one')      0.736135
('bar', 'two')     -0.859631
('baz', 'one')     -0.424100
('baz', 'two')     -0.776114
('foo', 'one')      1.279293
('foo', 'two')      0.943798
('qux', 'one')     -1.001859
('qux', 'two')      0.306546
```

The reason that the `MultiIndex` matters is that it can allow you to do grouping, selection, and reshaping operations as we will describe below and in subsequent areas of the documentation. As you will see in later sections, you can find

yourself working with hierarchically-indexed data without creating a `MultiIndex` explicitly yourself. However, when loading data from a file, you may wish to generate your own `MultiIndex` when preparing the data set.

Note that how the index is displayed by be controlled using the `multi_sparse` option in `pandas.set_printoptions`:

```
In [667]: pd.set_printoptions(multi_sparse=False)
```

```
In [668]: df
```

```
Out[668]:
<class 'pandas.core.frame.DataFrame'>
Index: 3 entries, A to C
Data columns:
('bar', 'one')      3 non-null values
('bar', 'two')      3 non-null values
('baz', 'one')      3 non-null values
('baz', 'two')      3 non-null values
('foo', 'one')      3 non-null values
('foo', 'two')      3 non-null values
('qux', 'one')      3 non-null values
('qux', 'two')      3 non-null values
dtypes: float64(8)
```

```
In [669]: pd.set_printoptions(multi_sparse=True)
```

## 7.4.2 Reconstructing the level labels

The method `get_level_values` will return a vector of the labels for each location at a particular level:

```
In [670]: index.get_level_values(0)
```

```
Out[670]: array([bar, bar, baz, baz, foo, foo, qux, qux], dtype=object)
```

```
In [671]: index.get_level_values('second')
```

```
Out[671]: array([one, two, one, two, one, two, one, two], dtype=object)
```

## 7.4.3 Basic indexing on axis with `MultiIndex`

One of the important features of hierarchical indexing is that you can select data by a “partial” label identifying a subgroup in the data. **Partial** selection “drops” levels of the hierarchical index in the result in a completely analogous way to selecting a column in a regular `DataFrame`:

```
In [672]: df['bar']
```

```
Out[672]:
second      one      two
A      -1.296337  0.150680
B       1.469725  1.304124
C      -0.938794  0.669142
```

```
In [673]: df['bar', 'one']
```

```
Out[673]:
A      -1.296337
B       1.469725
C      -0.938794
Name: ('bar', 'one')
```

```
In [674]: df['bar']['one']
```

```
Out [674]:
A    -1.296337
B     1.469725
C    -0.938794
Name: one
```

```
In [675]: s['qux']
Out [675]:
one    0.912560
two   -1.003401
```

## 7.4.4 Data alignment and using `reindex`

Operations between differently-indexed objects having `MultiIndex` on the axes will work as you expect; data alignment will work the same as an `Index` of tuples:

```
In [676]: s + s[:-2]
Out [676]:
bar one    -1.761968
    two     1.994577
baz one    -3.386631
    two    -0.358257
foo one    -3.196125
    two     1.873828
qux one         NaN
    two         NaN
```

```
In [677]: s + s[:,2]
Out [677]:
bar one    -1.761968
    two         NaN
baz one    -3.386631
    two         NaN
foo one    -3.196125
    two         NaN
qux one     1.825119
    two         NaN
```

`reindex` can be called with another `MultiIndex` or even a list or array of tuples:

```
In [678]: s.reindex(index[:3])
Out [678]:
first second
bar one    -0.880984
    two     0.997289
baz one    -1.693316
```

```
In [679]: s.reindex([('foo', 'two'), ('bar', 'one'), ('qux', 'one'), ('baz', 'one')])
Out [679]:
foo two     0.936914
bar one    -0.880984
qux one     0.912560
baz one    -1.693316
```

### 7.4.5 Advanced indexing with hierarchical index

Syntactically integrating `MultiIndex` in advanced indexing with `.ix` is a bit challenging, but we've made every effort to do so. for example the following works as you would expect:

```
In [680]: df = df.T
```

```
In [681]: df
```

```
Out[681]:
```

		A	B	C
first	second			
bar	one	-1.296337	1.469725	-0.938794
	two	0.150680	1.304124	0.669142
baz	one	0.123836	1.449735	-0.433567
	two	0.571764	0.203109	-0.273610
foo	one	1.555563	-1.032011	0.680433
	two	-0.823761	0.969818	-0.308450
qux	one	0.535420	-0.962723	-0.276099
	two	-1.032853	1.382083	-1.821168

```
In [682]: df.ix['bar']
```

```
Out[682]:
```

	A	B	C
second			
one	-1.296337	1.469725	-0.938794
two	0.150680	1.304124	0.669142

```
In [683]: df.ix['bar', 'two']
```

```
Out[683]:
```

A	0.150680
B	1.304124
C	0.669142

Name: ('bar', 'two')

“Partial” slicing also works quite nicely:

```
In [684]: df.ix['baz':'foo']
```

```
Out[684]:
```

		A	B	C
first	second			
baz	one	0.123836	1.449735	-0.433567
	two	0.571764	0.203109	-0.273610
foo	one	1.555563	-1.032011	0.680433
	two	-0.823761	0.969818	-0.308450

```
In [685]: df.ix[('baz', 'two'):(('qux', 'one'))]
```

```
Out[685]:
```

		A	B	C
first	second			
baz	two	0.571764	0.203109	-0.273610
foo	one	1.555563	-1.032011	0.680433
	two	-0.823761	0.969818	-0.308450
qux	one	0.535420	-0.962723	-0.276099

```
In [686]: df.ix[('baz', 'two'):'foo']
```

```
Out[686]:
```

		A	B	C
first	second			
baz	two	0.571764	0.203109	-0.273610

```
foo    one    1.555563 -1.032011  0.680433
      two    -0.823761  0.969818 -0.308450
```

Passing a list of labels or tuples works similar to reindexing:

```
In [687]: df.ix[(['bar', 'two'), ('qux', 'one')]]
```

```
Out[687]:
```

		A	B	C
first	second			
bar	two	0.15068	1.304124	0.669142
qux	one	0.53542	-0.962723	-0.276099

The following does not work, and it's not clear if it should or not:

```
>>> df.ix[['bar', 'qux']]
```

The code for implementing `.ix` makes every attempt to “do the right thing” but as you use it you may uncover corner cases or unintuitive behavior. If you do find something like this, do not hesitate to report the issue or ask on the mailing list.

## 7.4.6 Cross-section with hierarchical index

The `xs` method of `DataFrame` additionally takes a `level` argument to make selecting data at a particular level of a `MultiIndex` easier.

```
In [688]: df.xs('one', level='second')
```

```
Out[688]:
```

		A	B	C
first				
bar		-1.296337	1.469725	-0.938794
baz		0.123836	1.449735	-0.433567
foo		1.555563	-1.032011	0.680433
qux		0.535420	-0.962723	-0.276099

## 7.4.7 Advanced reindexing and alignment with hierarchical index

The parameter `level` has been added to the `reindex` and `align` methods of pandas objects. This is useful to broadcast values across a level. For instance:

```
In [689]: midx = MultiIndex(levels=[['zero', 'one'], ['x', 'y']],
.....:                      labels=[[1, 1, 0, 0], [1, 0, 1, 0]])
.....:
```

```
In [690]: df = DataFrame(randn(4, 2), index=midx)
```

```
In [691]: print df
```

		0	1
one	y	0.307453	-0.906534
	x	-1.505397	1.392009
zero	y	-0.027793	-0.631023
	x	-0.662357	2.725042

```
In [692]: df2 = df.mean(level=0)
```

```
In [693]: print df2
```

		0	1
--	--	---	---

```
zero -0.345075  1.047010
one  -0.598972  0.242737
```

```
In [694]: print df2.reindex(df.index, level=0)
           0      1
one  y -0.598972  0.242737
     x -0.598972  0.242737
zero y -0.345075  1.047010
     x -0.345075  1.047010
```

```
In [695]: df_aligned, df2_aligned = df.align(df2, level=0)
```

```
In [696]: print df_aligned
           0      1
one  y  0.307453 -0.906534
     x -1.505397  1.392009
zero y -0.027793 -0.631023
     x -0.662357  2.725042
```

```
In [697]: print df2_aligned
           0      1
one  y -0.598972  0.242737
     x -0.598972  0.242737
zero y -0.345075  1.047010
     x -0.345075  1.047010
```

## 7.4.8 The need for sortedness

**Caveat emptor:** the present implementation of `MultiIndex` requires that the labels be sorted for some of the slicing / indexing routines to work correctly. You can think about breaking the axis into unique groups, where at the hierarchical level of interest, each distinct group shares a label, but no two have the same label. However, the `MultiIndex` does not enforce this: **you are responsible for ensuring that things are properly sorted**. There is an important new method `sortlevel` to sort an axis within a `MultiIndex` so that its labels are grouped and sorted by the original ordering of the associated factor at that level. Note that this does not necessarily mean the labels will be sorted lexicographically!

```
In [698]: import random; random.shuffle(tuples)
```

```
In [699]: s = Series(randn(8), index=MultiIndex.from_tuples(tuples))
```

```
In [700]: s
Out[700]:
foo two    -1.847240
qux two    -0.529247
bar one     0.614656
qux one    -1.590742
baz one    -0.156479
bar two    -1.696377
baz two     0.819712
foo one    -2.107728
```

```
In [701]: s.sortlevel(0)
Out[701]:
bar one     0.614656
      two    -1.696377
baz one    -0.156479
```

```
      two      0.819712
foo one  -2.107728
      two  -1.847240
qux one  -1.590742
      two  -0.529247
```

```
In [702]: s.sortlevel(1)
```

```
Out[702]:
bar one      0.614656
baz one     -0.156479
foo one     -2.107728
qux one     -1.590742
bar two     -1.696377
baz two      0.819712
foo two     -1.847240
qux two     -0.529247
```

Note, you may also pass a level name to `sortlevel` if the MultiIndex levels are named.

```
In [703]: s.index.names = ['L1', 'L2']
```

```
In [704]: s.sortlevel(level='L1')
```

```
Out[704]:
L1  L2
bar one      0.614656
      two     -1.696377
baz one     -0.156479
      two      0.819712
foo one     -2.107728
      two     -1.847240
qux one     -1.590742
      two     -0.529247
```

```
In [705]: s.sortlevel(level='L2')
```

```
Out[705]:
L1  L2
bar one      0.614656
baz one     -0.156479
foo one     -2.107728
qux one     -1.590742
bar two     -1.696377
baz two      0.819712
foo two     -1.847240
qux two     -0.529247
```

Some indexing will work even if the data are not sorted, but will be rather inefficient and will also return a copy of the data rather than a view:

```
In [706]: s['qux']
```

```
Out[706]:
L2
two  -0.529247
one  -1.590742
```

```
In [707]: s.sortlevel(1)['qux']
```

```
Out[707]:
L2
one  -1.590742
two  -0.529247
```

On higher dimensional objects, you can sort any of the other axes by level if they have a MultiIndex:

```
In [708]: df.T.sortlevel(1, axis=1)
Out[708]:
```

	zero	one	zero	one
	x	x	y	y
0	-0.662357	-1.505397	-0.027793	0.307453
1	2.725042	1.392009	-0.631023	-0.906534

The MultiIndex object has code to **explicitly check the sort depth**. Thus, if you try to index at a depth at which the index is not sorted, it will raise an exception. Here is a concrete example to illustrate this:

```
In [709]: tuples = [('a', 'a'), ('a', 'b'), ('b', 'a'), ('b', 'b')]
```

```
In [710]: idx = MultiIndex.from_tuples(tuples)
```

```
In [711]: idx.lexsort_depth
Out[711]: 2
```

```
In [712]: reordered = idx[[1, 0, 3, 2]]
```

```
In [713]: reordered.lexsort_depth
Out[713]: 1
```

```
In [714]: s = Series(randn(4), index=reordered)
```

```
In [715]: s.ix['a':'a']
Out[715]:
```

a	b	-0.488326
a	a	0.851918

However:

```
>>> s.ix[('a', 'b'):( 'b', 'a')]
Exception: MultiIndex lexsort depth 1, key was length 2
```

## 7.4.9 Swapping levels with `swaplevel`

The `swaplevel` function can switch the order of two levels:

```
In [716]: df[:5]
Out[716]:
```

		0	1
one	y	0.307453	-0.906534
	x	-1.505397	1.392009
zero	y	-0.027793	-0.631023
	x	-0.662357	2.725042

```
In [717]: df[:5].swaplevel(0, 1, axis=0)
Out[717]:
```

		0	1
y	one	0.307453	-0.906534
x	one	-1.505397	1.392009
y	zero	-0.027793	-0.631023
x	zero	-0.662357	2.725042



### 7.4.10 Reordering levels with `reorder_levels`

The `reorder_levels` function generalizes the `swaplevel` function, allowing you to permute the hierarchical index levels in one step:

```
In [718]: df[:5].reorder_levels([1,0], axis=0)
Out[718]:
```

	0	1
y one	0.307453	-0.906534
x one	-1.505397	1.392009
y zero	-0.027793	-0.631023
x zero	-0.662357	2.725042

### 7.4.11 Some gory internal details

Internally, the `MultiIndex` consists of a few things: the **levels**, the integer **labels**, and the level **names**:

```
In [719]: index
Out[719]:
MultiIndex
[('bar', 'one') ('bar', 'two') ('baz', 'one') ('baz', 'two')
 ('foo', 'one') ('foo', 'two') ('qux', 'one') ('qux', 'two')]

In [720]: index.levels
Out[720]: [Index([bar, baz, foo, qux], dtype=object), Index([one, two], dtype=object)]

In [721]: index.labels
Out[721]: [array([0, 0, 1, 1, 2, 2, 3, 3]), array([0, 1, 0, 1, 0, 1, 0, 1])]

In [722]: index.names
Out[722]: ['first', 'second']
```

You can probably guess that the labels determine which unique element is identified with that location at each layer of the index. It's important to note that sortedness is determined **solely** from the integer labels and does not check (or care) whether the levels themselves are sorted. Fortunately, the constructors `from_tuples` and `from_arrays` ensure that this is true, but if you compute the levels and labels yourself, please be careful.

## 7.5 Adding an index to an existing DataFrame

Occasionally you will load or create a data set into a `DataFrame` and want to add an index after you've already done so. There are a couple of different ways.

### 7.5.1 Add an index using `DataFrame` columns

`DataFrame` has a `set_index` method which takes a column name (for a regular `Index`) or a list of column names (for a `MultiIndex`), to create a new, indexed `DataFrame`:

```
In [723]: data
Out[723]:
```

	a	b	c	d
0	bar	one	z	1
1	bar	two	y	2
2	foo	one	x	3
3	foo	two	w	4

```
In [724]: indexed1 = data.set_index('c')
```

```
In [725]: indexed1
```

```
Out[725]:
```

	a	b	d
c			
z	bar	one	1
y	bar	two	2
x	foo	one	3
w	foo	two	4

```
In [726]: indexed2 = data.set_index(['a', 'b'])
```

```
In [727]: indexed2
```

```
Out[727]:
```

		c	d
a	b		
bar	one	z	1
	two	y	2
foo	one	x	3
	two	w	4

The append keyword option allow you to keep the existing index and append the given columns to a MultiIndex:

```
In [728]: frame = data.set_index('c', drop=False)
```

```
In [729]: frame = frame.set_index(['a', 'b'], append=True)
```

```
In [730]: frame
```

```
Out[730]:
```

			c	d
c	a	b		
z	bar	one	z	1
y	bar	two	y	2
x	foo	one	x	3
w	foo	two	w	4

Other options in `set_index` allow you not drop the index columns or to add the index in-place (without creating a new object):

```
In [731]: data.set_index('c', drop=False)
```

```
Out[731]:
```

	a	b	c	d
c				
z	bar	one	z	1
y	bar	two	y	2
x	foo	one	x	3
w	foo	two	w	4

```
In [732]: df = data.set_index(['a', 'b'], inplace=True)
```

```
In [733]: data
```

```
Out[733]:
```

		c	d
a	b		
bar	one	z	1
	two	y	2
foo	one	x	3

```
two w 4
```

## 7.5.2 Remove / reset the index, `reset_index`

As a convenience, there is a new function on `DataFrame` called `reset_index` which transfers the index values into the `DataFrame`'s columns and sets a simple integer index. This is the inverse operation to `set_index`

```
In [734]: df
```

```
Out[734]:
```

		c	d
a	b		
bar	one	z	1
	two	y	2
foo	one	x	3
	two	w	4

```
In [735]: df.reset_index()
```

```
Out[735]:
```

	a	b	c	d
0	bar	one	z	1
1	bar	two	y	2
2	foo	one	x	3
3	foo	two	w	4

The output is more similar to a SQL table or a record array. The names for the columns derived from the index are the ones stored in the `names` attribute.

You can use the `level` keyword to remove only a portion of the index:

```
In [736]: frame
```

```
Out[736]:
```

		c	d
c	a	b	
z	bar	one	z 1
y	bar	two	y 2
x	foo	one	x 3
w	foo	two	w 4

```
In [737]: frame.reset_index(level=1)
```

```
Out[737]:
```

		a	c	d
c	b			
z	one	bar	z	1
y	two	bar	y	2
x	one	foo	x	3
w	two	foo	w	4

`reset_index` takes an optional parameter `drop` which if true simply discards the index, instead of putting index values in the `DataFrame`'s columns.

---

**Note:** The `reset_index` method used to be called `delevel` which is now deprecated.

---

## 7.5.3 Adding an ad hoc index

If you create an index yourself, you can just assign it to the `index` field:

```
df.index = index
```

## 7.6 Indexing internal details

---

**Note:** The following is largely relevant for those actually working on the pandas codebase. And the source code is still the best place to look at the specifics of how things are implemented.

---

In pandas there are a few objects implemented which can serve as valid containers for the axis labels:

- `Index`: the generic “ordered set” object, an ndarray of object dtype assuming nothing about its contents. The labels must be hashable (and likely immutable) and unique. Populates a dict of label to location in Cython to do  $O(1)$  lookups.
- `Int64Index`: a version of `Index` highly optimized for 64-bit integer data, such as time stamps
- `MultiIndex`: the standard hierarchical index object
- `date_range`: fixed frequency date range generated from a time rule or `DateOffset`. An ndarray of Python datetime objects

The motivation for having an `Index` class in the first place was to enable different implementations of indexing. This means that it’s possible for you, the user, to implement a custom `Index` subclass that may be better suited to a particular application than the ones provided in pandas. For example, we plan to add a more efficient datetime index which leverages the new `numpy.datetime64` dtype in the relatively near future.

From an internal implementation point of view, the relevant methods that an `Index` must define are one or more of the following (depending on how incompatible the new object internals are with the `Index` functions):

- `get_loc`: returns an “indexer” (an integer, or in some cases a slice object) for a label
- `slice_locs`: returns the “range” to slice between two labels
- `get_indexer`: Computes the indexing vector for reindexing / data alignment purposes. See the source / docstrings for more on this
- `reindex`: Does any pre-conversion of the input index then calls `get_indexer`
- `union, intersection`: computes the union or intersection of two `Index` objects
- `insert`: Inserts a new label into an `Index`, yielding a new object
- `delete`: Delete a label, yielding a new object
- `drop`: Deletes a set of labels
- `take`: Analogous to `ndarray.take`

# COMPUTATIONAL TOOLS

## 8.1 Statistical functions

### 8.1.1 Percent Change

Both `Series` and `DataFrame` has a method `pct_change` to compute the percent change over a given number of periods (using `fill_method` to fill NA/null values).

```
In [187]: ser = Series(randn(8))
```

```
In [188]: ser.pct_change()
```

```
Out[188]:
0          NaN
1   -1.602976
2    4.334938
3   -0.247456
4   -2.067345
5   -1.142903
6   -1.688214
7   -9.759729
```

```
In [189]: df = DataFrame(randn(10, 4))
```

```
In [190]: df.pct_change(periods=3)
```

```
Out[190]:
      0         1         2         3
0     NaN     NaN     NaN     NaN
1     NaN     NaN     NaN     NaN
2     NaN     NaN     NaN     NaN
3 -0.218320 -1.054001  1.987147 -0.510183
4 -0.439121 -1.816454  0.649715 -4.822809
5 -0.127833 -3.042065 -5.866604 -1.776977
6 -2.596833 -1.959538 -2.111697 -3.798900
7 -0.117826 -2.169058  0.036094 -0.067696
8  2.492606 -1.357320 -1.205802 -1.558697
9 -1.012977  2.324558 -1.003744 -0.371806
```

### 8.1.2 Covariance

The `Series` object has a method `cov` to compute covariance between series (excluding NA/null values).

```
In [191]: s1 = Series(randn(1000))
```

```
In [192]: s2 = Series(randn(1000))
```

```
In [193]: s1.cov(s2)
```

```
Out[193]: 0.00068010881743109321
```

Analogously, `DataFrame` has a method `cov` to compute pairwise covariances among the series in the `DataFrame`, also excluding NA/null values.

```
In [194]: frame = DataFrame(randn(1000, 5), columns=['a', 'b', 'c', 'd', 'e'])
```

```
In [195]: frame.cov()
```

```
Out[195]:
```

	a	b	c	d	e
a	1.000882	-0.003177	-0.002698	-0.006889	0.031912
b	-0.003177	1.024721	0.000191	0.009212	0.000857
c	-0.002698	0.000191	0.950735	-0.031743	-0.005087
d	-0.006889	0.009212	-0.031743	1.002983	-0.047952
e	0.031912	0.000857	-0.005087	-0.047952	1.042487

## 8.1.3 Correlation

Several methods for computing correlations are provided. Several kinds of correlation methods are provided:

Method name	Description
pearson (default)	Standard correlation coefficient
kendall	Kendall Tau correlation coefficient
spearman	Spearman rank correlation coefficient

All of these are currently computed using pairwise complete observations.

```
In [196]: frame = DataFrame(randn(1000, 5), columns=['a', 'b', 'c', 'd', 'e'])
```

```
In [197]: frame.ix[:,2] = np.nan
```

```
# Series with Series
```

```
In [198]: frame['a'].corr(frame['b'])
```

```
Out[198]: 0.010052135416653445
```

```
In [199]: frame['a'].corr(frame['b'], method='spearman')
```

```
Out[199]: -0.0097383749534998149
```

```
# Pairwise correlation of DataFrame columns
```

```
In [200]: frame.corr()
```

```
Out[200]:
```

	a	b	c	d	e
a	1.000000	0.010052	-0.047750	-0.031461	-0.025285
b	0.010052	1.000000	-0.014172	-0.020590	-0.001930
c	-0.047750	-0.014172	1.000000	0.006373	-0.049479
d	-0.031461	-0.020590	0.006373	1.000000	-0.012379
e	-0.025285	-0.001930	-0.049479	-0.012379	1.000000

Note that non-numeric columns will be automatically excluded from the correlation calculation.

A related method `corrwith` is implemented on `DataFrame` to compute the correlation between like-labeled `Series` contained in different `DataFrame` objects.

```

In [201]: index = ['a', 'b', 'c', 'd', 'e']

In [202]: columns = ['one', 'two', 'three', 'four']

In [203]: df1 = DataFrame(randn(5, 4), index=index, columns=columns)

In [204]: df2 = DataFrame(randn(4, 4), index=index[:4], columns=columns)

In [205]: df1.corrwith(df2)
Out[205]:
one      0.803464
two      0.142469
three   -0.498774
four     0.806420

In [206]: df2.corrwith(df1, axis=1)
Out[206]:
a      0.011572
b      0.388066
c     -0.335819
d      0.232412
e           NaN

```

### 8.1.4 Data ranking

The rank method produces a data ranking with ties being assigned the mean of the ranks (by default) for the group:

```

In [207]: s = Series(np.random.randn(5), index=list('abcde'))

In [208]: s['d'] = s['b'] # so there's a tie

In [209]: s.rank()
Out[209]:
a      2.0
b      4.5
c      3.0
d      4.5
e      1.0

```

rank is also a DataFrame method and can rank either the rows (axis=0) or the columns (axis=1). NaN values are excluded from the ranking.

```

In [210]: df = DataFrame(np.random.randn(10, 6))

In [211]: df[4] = df[2][:5] # some ties

In [212]: df
Out[212]:
   0      1      2      3      4      5
0  0.085011 -0.459422 -1.660917 -1.913019 -1.660917  0.833479
1 -0.557052  0.775425  0.003794  0.555351  0.003794 -1.169977
2  0.815695 -0.295737 -0.534290  0.068917 -0.534290 -0.513855
3  1.465947  0.021757  0.523224 -0.439297  0.523224 -0.959568
4 -0.678378  0.091855  1.337956  0.792551  1.337956  0.711776
5 -0.190285  0.187520 -0.355562  1.730964      NaN -1.362312
6 -0.776678 -2.082637 -0.165877  0.357163      NaN  0.631662
7 -1.295037  0.367656 -1.886797 -0.531790      NaN  1.270408

```

```
8  1.106052  0.848312 -0.613544  1.338296      NaN -1.150652
9  0.309979  1.088439  0.920366 -0.750322      NaN  1.563956
```

```
In [213]: df.rank(1)
```

```
Out[213]:
   0  1  2  3  4  5
0  5  4  2.5  1  2.5  6
1  2  6  3.5  5  3.5  1
2  6  4  1.5  5  1.5  3
3  6  3  4.5  2  4.5  1
4  1  2  5.5  4  5.5  3
5  3  4  2.0  5  NaN  1
6  2  1  3.0  4  NaN  5
7  2  4  1.0  3  NaN  5
8  4  3  2.0  5  NaN  1
9  2  4  3.0  1  NaN  5
```

`rank` optionally takes a parameter `ascending` which by default is `true`; when `false`, data is reverse-ranked, with larger values assigned a smaller rank.

`rank` supports different tie-breaking methods, specified with the `method` parameter:

- `average` : average rank of tied group
- `min` : lowest rank in the group
- `max` : highest rank in the group
- `first` : ranks assigned in the order they appear in the array

---

**Note:** These methods are significantly faster (around 10-20x) than `scipy.stats.rankdata`.

---

## 8.2 Moving (rolling) statistics / moments

For working with time series data, a number of functions are provided for computing common *moving* or *rolling* statistics. Among these are count, sum, mean, median, correlation, variance, covariance, standard deviation, skewness, and kurtosis. All of these methods are in the `pandas` namespace, but otherwise they can be found in `pandas.stats.moments`.

Function	Description
<code>rolling_count</code>	Number of non-null observations
<code>rolling_sum</code>	Sum of values
<code>rolling_mean</code>	Mean of values
<code>rolling_median</code>	Arithmetic median of values
<code>rolling_min</code>	Minimum
<code>rolling_max</code>	Maximum
<code>rolling_std</code>	Unbiased standard deviation
<code>rolling_var</code>	Unbiased variance
<code>rolling_skew</code>	Unbiased skewness (3rd moment)
<code>rolling_kurt</code>	Unbiased kurtosis (4th moment)
<code>rolling_quantile</code>	Sample quantile (value at %)
<code>rolling_apply</code>	Generic apply
<code>rolling_cov</code>	Unbiased covariance (binary)
<code>rolling_corr</code>	Correlation (binary)
<code>rolling_corr_pairwise</code>	Pairwise correlation of DataFrame columns



Generally these methods all have the same interface. The binary operators (e.g. `rolling_corr`) take two Series or DataFrames. Otherwise, they all accept the following arguments:

- `window`: size of moving window
- `min_periods`: threshold of non-null data points to require (otherwise result is NA)
- `freq`: optionally specify a `ref`: *frequency string* `<timeseries.alias>` or *DateOffset* to pre-conform the data to. Note that prior to pandas v0.8.0, a keyword argument `time_rule` was used instead of `freq` that referred to the legacy time rule constants

These functions can be applied to ndarrays or Series objects:

```
In [214]: ts = Series(randn(1000), index=date_range('1/1/2000', periods=1000))
```

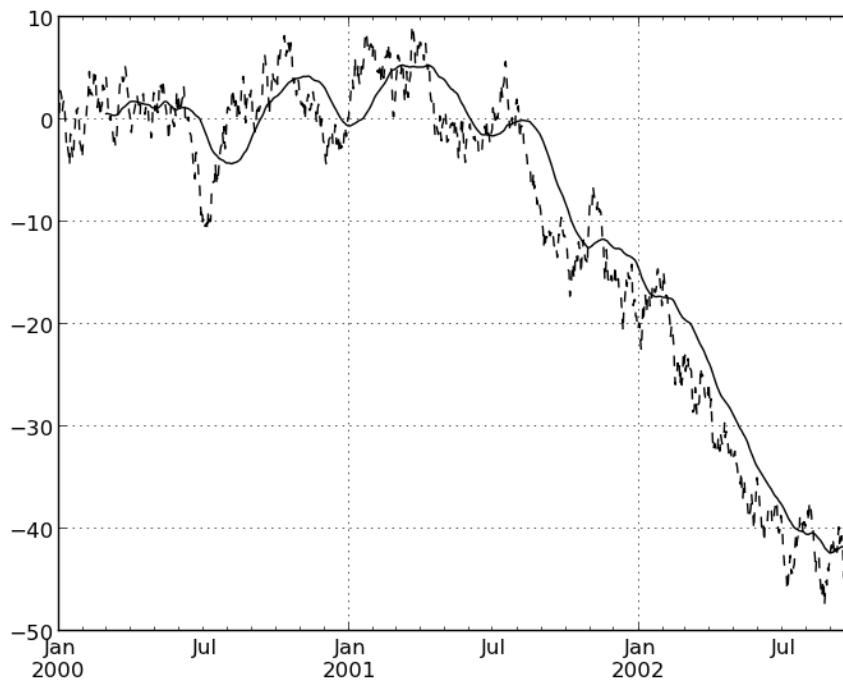
```
In [215]: ts = ts.cumsum()
```

```
In [216]: ts.plot(style='k--')
```

```
Out[216]: <matplotlib.axes.AxesSubplot at 0x5c4e610>
```

```
In [217]: rolling_mean(ts, 60).plot(style='k')
```

```
Out[217]: <matplotlib.axes.AxesSubplot at 0x5c4e610>
```



They can also be applied to DataFrame objects. This is really just syntactic sugar for applying the moving window operator to all of the DataFrame's columns:

```
In [218]: df = DataFrame(randn(1000, 4), index=ts.index,
.....:                  columns=['A', 'B', 'C', 'D'])
.....:
```

```
In [219]: df = df.cumsum()
```

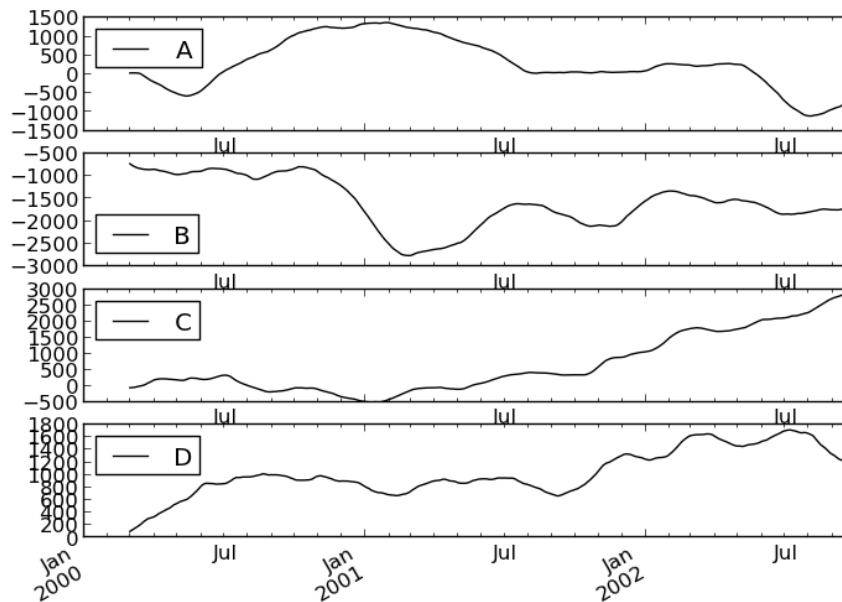
```
In [220]: rolling_sum(df, 60).plot(subplots=True)
```

```
Out[220]:
array([Axes(0.125,0.747826;0.775x0.152174),
```

```

Axes(0.125,0.565217;0.775x0.152174),
Axes(0.125,0.382609;0.775x0.152174)], dtype=object)

```

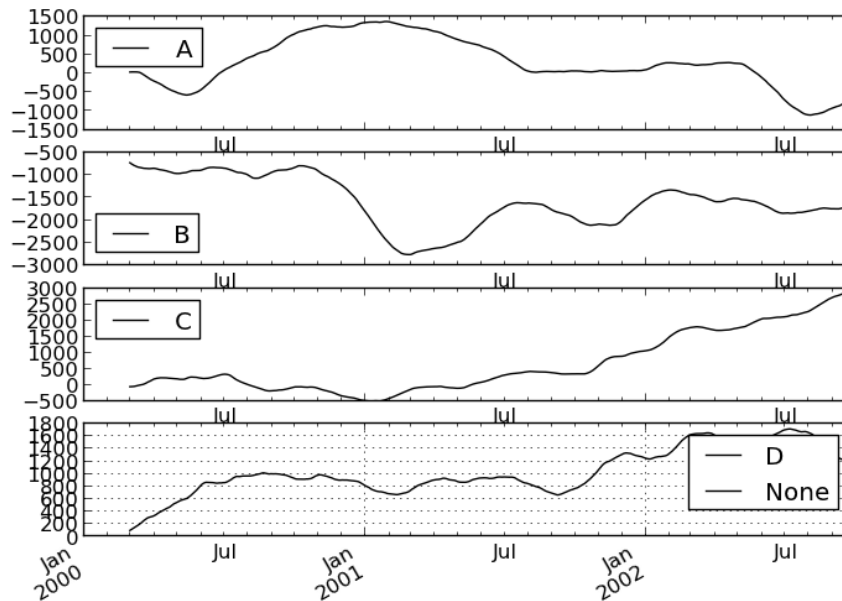


The `rolling_apply` function takes an extra `func` argument and performs generic rolling computations. The `func` argument should be a single function that produces a single value from an ndarray input. Suppose we wanted to compute the mean absolute deviation on a rolling basis:

```
In [221]: mad = lambda x: np.fabs(x - x.mean()).mean()
```

```
In [222]: rolling_apply(ts, 60, mad).plot(style='k')
```

```
Out[222]: <matplotlib.axes.AxesSubplot at 0x6021810>
```



## 8.2.1 Binary rolling moments

`rolling_cov` and `rolling_corr` can compute moving window statistics about two `Series` or any combination of `DataFrame/Series` or `DataFrame/DataFrame`. Here is the behavior in each case:

- two `Series`: compute the statistic for the pairing
- `DataFrame/Series`: compute the statistics for each column of the `DataFrame` with the passed `Series`, thus returning a `DataFrame`
- `DataFrame/DataFrame`: compute statistic for matching column names, returning a `DataFrame`

For example:

```
In [223]: df2 = df[:20]
```

```
In [224]: rolling_corr(df2, df2['B'], window=5)
```

```
Out[224]:
```

	A	B	C	D
2000-01-01	NaN	NaN	NaN	NaN
2000-01-02	NaN	NaN	NaN	NaN
2000-01-03	NaN	NaN	NaN	NaN
2000-01-04	NaN	NaN	NaN	NaN
2000-01-05	0.703188	1	-0.746130	0.714265
2000-01-06	0.065322	1	-0.209789	0.635360
2000-01-07	-0.429914	1	-0.100807	0.266005
2000-01-08	-0.387498	1	0.512321	0.592033
2000-01-09	0.442207	1	0.570186	-0.653242
2000-01-10	0.572983	1	0.713876	-0.366806
2000-01-11	0.325889	1	0.899489	-0.337436
2000-01-12	-0.389584	1	0.482351	0.246871
2000-01-13	-0.714206	1	-0.593838	0.090279
2000-01-14	-0.933238	1	-0.936087	0.471866
2000-01-15	-0.991959	1	-0.943218	0.637434
2000-01-16	-0.645081	1	-0.520788	0.322264
2000-01-17	-0.348338	1	-0.183528	0.385915
2000-01-18	0.193914	1	-0.308346	-0.157765
2000-01-19	0.465424	1	-0.072219	-0.714273
2000-01-20	0.645630	1	0.211302	-0.651308

## 8.2.2 Computing rolling pairwise correlations

In financial data analysis and other fields it's common to compute correlation matrices for a collection of time series. More difficult is to compute a moving-window correlation matrix. This can be done using the `rolling_corr_pairwise` function, which yields a `Panel` whose items are the dates in question:

```
In [225]: correls = rolling_corr_pairwise(df, 50)
```

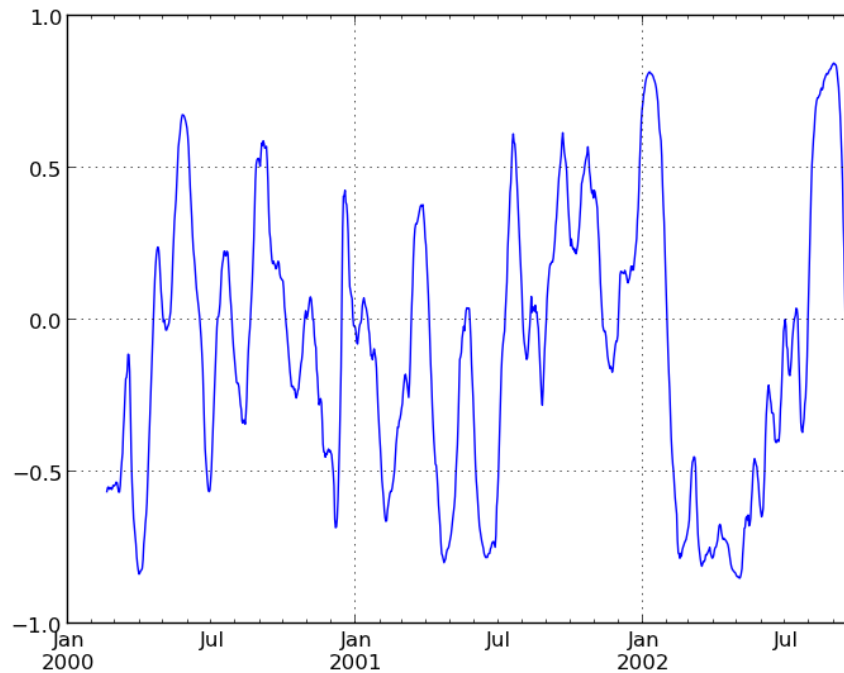
```
In [226]: correls[df.index[-50]]
```

```
Out[226]:
```

	A	B	C	D
A	1.000000	0.289597	0.673828	-0.589002
B	0.289597	1.000000	-0.041244	0.204692
C	0.673828	-0.041244	1.000000	-0.848632
D	-0.589002	0.204692	-0.848632	1.000000

You can efficiently retrieve the time series of correlations between two columns using `ix` indexing:

```
In [227]: correls.ix[:, 'A', 'C'].plot()
Out[227]: <matplotlib.axes.AxesSubplot at 0x62b23d0>
```



### 8.3 Exponentially weighted moment functions

A related set of functions are exponentially weighted versions of many of the above statistics. A number of EW (exponentially weighted) functions are provided using the blending method. For example, where  $y_t$  is the result and  $x_t$  the input, we compute an exponentially weighted moving average as

$$y_t = \alpha y_{t-1} + (1 - \alpha)x_t$$

One must have  $0 < \alpha \leq 1$ , but rather than pass  $\alpha$  directly, it's easier to think about either the **span** or **center of mass (com)** of an EW moment:

$$\alpha = \begin{cases} \frac{2}{s+1}, s = \text{span} \\ \frac{1}{c+1}, c = \text{center of mass} \end{cases}$$

You can pass one or the other to these functions but not both. **Span** corresponds to what is commonly called a “20-day EW moving average” for example. **Center of mass** has a more physical interpretation. For example, **span** = 20 corresponds to **com** = 9.5. Here is the list of functions available:

Function	Description
ewma	EW moving average
ewvar	EW moving variance
ewstd	EW moving standard deviation
ewmcorr	EW moving correlation
ewmcov	EW moving covariance

Here are an example for a univariate time series:

```
In [228]: plt.close('all')
In [229]: ts.plot(style='k--')
Out[229]: <matplotlib.axes.AxesSubplot at 0x6d9d090>
In [230]: ewma(ts, span=20).plot(style='k')
Out[230]: <matplotlib.axes.AxesSubplot at 0x6d9d090>
```



**Note:** The EW functions perform a standard adjustment to the initial observations whereby if there are fewer observations than called for in the span, those observations are reweighted accordingly.

## 8.4 Linear and panel regression

**Note:** We plan to move this functionality to [statsmodels](#) for the next release. Some of the result attributes may change names in order to foster naming consistency with the rest of statsmodels. We will provide every effort to provide compatibility with older versions of pandas, however.

We have implemented a very fast set of *moving-window linear regression* classes in pandas. Two different types of regressions are supported:

- Standard ordinary least squares (OLS) multiple regression
- Multiple regression (OLS-based) on [panel data](#) including with fixed-effects (also known as entity or individual effects) or time-effects.

Both kinds of linear models are accessed through the `ols` function in the pandas namespace. They all take the following arguments to specify either a static (full sample) or dynamic (moving window) regression:

- `window_type`: 'full sample' (default), 'expanding', or 'rolling'

- `window`: size of the moving window in the `window_type='rolling'` case. If `window` is specified, `window_type` will be automatically set to `'rolling'`
- `min_periods`: minimum number of time periods to require to compute the regression coefficients

Generally speaking, the `ols` works by being given a `y` (response) object and an `x` (predictors) object. These can take many forms:

- `y`: a `Series`, `ndarray`, or `DataFrame` (panel model)
- `x`: `Series`, `DataFrame`, `dict` of `Series`, `dict` of `DataFrame` or `Panel`

Based on the types of `y` and `x`, the model will be inferred to either a panel model or a regular linear model. If the `y` variable is a `DataFrame`, the result will be a panel model. In this case, the `x` variable must either be a `Panel`, or a `dict` of `DataFrame` (which will be coerced into a `Panel`).

### 8.4.1 Standard OLS regression

Let's pull in some sample data:

```
In [231]: from pandas.io.data import DataReader

In [232]: symbols = ['MSFT', 'GOOG', 'AAPL']

In [233]: data = dict((sym, DataReader(sym, "yahoo"))
.....:                 for sym in symbols)
.....:

In [234]: panel = Panel(data).swapaxes('items', 'minor')

In [235]: close_px = panel['Close']

# convert closing prices to returns
In [236]: rets = close_px / close_px.shift(1) - 1

In [237]: rets.info()
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 251 entries, 2011-07-25 00:00:00 to 2012-07-20 00:00:00
Data columns:
AAPL      250  non-null values
GOOG      250  non-null values
MSFT      250  non-null values
dtypes: float64(3)
```

Let's do a static regression of AAPL returns on GOOG returns:

```
In [238]: model = ols(y=rets['AAPL'], x=rets.ix[:, ['GOOG']])

In [239]: model
Out[239]:
-----Summary of Regression Analysis-----
Formula: Y ~ <GOOG> + <intercept>
Number of Observations:      250
Number of Degrees of Freedom: 2
R-squared:      0.3735
Adj R-squared:   0.3710
Rmse:           0.0148
F-stat (1, 248):  147.8494, p-value:      0.0000
Degrees of Freedom: model 1, resid 248
```

```
-----Summary of Estimated Coefficients-----
      Variable      Coef      Std Err      t-stat      p-value      CI 2.5%      CI 97.5%
-----
      GOOG      0.6594      0.0542      12.16      0.0000      0.5531      0.7657
      intercept      0.0018      0.0009      1.89      0.0594      -0.0001      0.0036
-----End of Summary-----
```

```
In [240]: model.beta
Out[240]:
GOOG      0.659439
intercept 0.001775
```

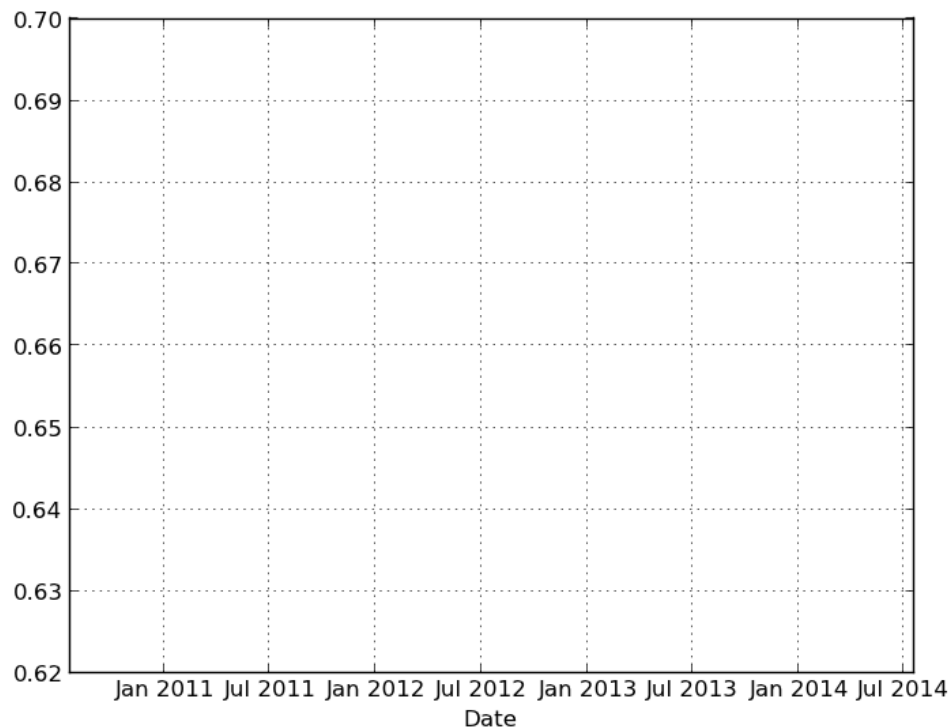
If we had passed a Series instead of a DataFrame with the single GOOG column, the model would have assigned the generic name `x` to the sole right-hand side variable.

We can do a moving window regression to see how the relationship changes over time:

```
In [241]: model = ols(y=rets['AAPL'], x=rets.ix[:, ['GOOG']],
.....:               window=250)
.....:
```

```
# just plot the coefficient for GOOG
```

```
In [242]: model.beta['GOOG'].plot()
Out[242]: <matplotlib.axes.AxesSubplot at 0x7110410>
```



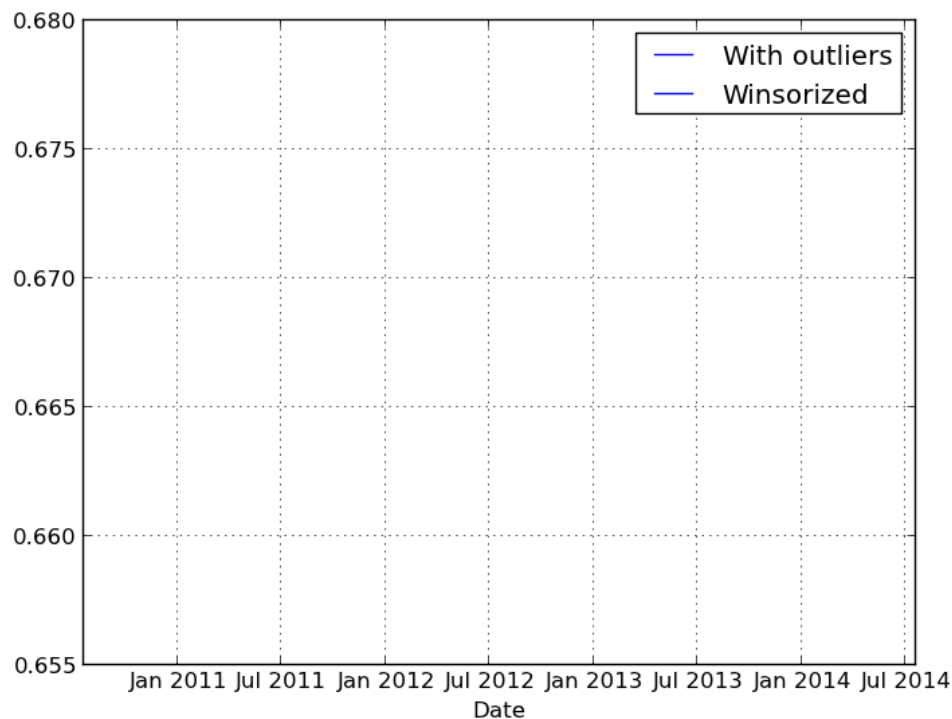
It looks like there are some outliers rolling in and out of the window in the above regression, influencing the results. We could perform a simple [winsorization](#) at the 3 STD level to trim the impact of outliers:

```
In [243]: winz = rets.copy()

In [244]: std_1year = rolling_std(rets, 250, min_periods=20)

# cap at 3 * 1 year standard deviation
```

```
In [245]: cap_level = 3 * np.sign(winz) * std_1year
In [246]: winz[np.abs(winz) > 3 * std_1year] = cap_level
In [247]: winz_model = ols(y=winz['AAPL'], x=winz.ix[:, ['GOOG']],
.....:                    window=250)
.....:
In [248]: model.beta['GOOG'].plot(label="With outliers")
Out[248]: <matplotlib.axes.AxesSubplot at 0x70ff850>
In [249]: winz_model.beta['GOOG'].plot(label="Winsorized"); plt.legend(loc='best')
Out[249]: <matplotlib.legend.Legend at 0x7d392d0>
```



So in this simple example we see the impact of winsorization is actually quite significant. Note the correlation after winsorization remains high:

```
In [250]: winz.corrwith(rets)
Out[250]:
AAPL    0.995128
GOOG    0.997097
MSFT    0.999961
```

Multiple regressions can be run by passing a DataFrame with multiple columns for the predictors x:

```
In [251]: ols(y=winz['AAPL'], x=winz.drop(['AAPL'], axis=1))
Out[251]:
-----Summary of Regression Analysis-----
Formula: Y ~ <GOOG> + <MSFT> + <intercept>
Number of Observations:      250
Number of Degrees of Freedom:  3
R-squared:                    0.4389
Adj R-squared:                0.4343
```



```

Rmse:                0.0137
F-stat (2, 247):      96.5973, p-value:      0.0000
Degrees of Freedom: model 2, resid 247
-----Summary of Estimated Coefficients-----
      Variable      Coef      Std Err      t-stat      p-value      CI 2.5%      CI 97.5%
-----
          GOOG      0.5411      0.0609       8.88      0.0000       0.4217       0.6605
          MSFT      0.2702      0.0655       4.13      0.0001       0.1418       0.3985
      intercept      0.0015      0.0009       1.73      0.0857      -0.0002       0.0032
-----End of Summary-----

```

## 8.4.2 Panel regression

We've implemented moving window panel regression on potentially unbalanced panel data (see [this article](#) if this means nothing to you). Suppose we wanted to model the relationship between the magnitude of the daily return and trading volume among a group of stocks, and we want to pool all the data together to run one big regression. This is actually quite easy:

```

# make the units somewhat comparable
In [252]: volume = panel['Volume'] / 1e8

In [253]: model = ols(y=volume, x={'return' : np.abs(rets)})

In [254]: model
Out[254]:
-----Summary of Regression Analysis-----
Formula: Y ~ <return> + <intercept>
Number of Observations:      750
Number of Degrees of Freedom:  2
R-squared:                    0.0306
Adj R-squared:                0.0293
Rmse:                        0.2414
F-stat (1, 748):              23.6403, p-value:      0.0000
Degrees of Freedom: model 1, resid 748
-----Summary of Estimated Coefficients-----
      Variable      Coef      Std Err      t-stat      p-value      CI 2.5%      CI 97.5%
-----
      return      3.6646      0.7537       4.86      0.0000       2.1873       5.1418
      intercept      0.2043      0.0130      15.66      0.0000       0.1788       0.2299
-----End of Summary-----

```

In a panel model, we can insert dummy (0-1) variables for the “entities” involved (here, each of the stocks) to account the a entity-specific effect (intercept):

```

In [255]: fe_model = ols(y=volume, x={'return' : np.abs(rets)},
.....:                  entity_effects=True)
.....:

In [256]: fe_model
Out[256]:
-----Summary of Regression Analysis-----
Formula: Y ~ <return> + <FE_GOOG> + <FE_MSFT> + <intercept>
Number of Observations:      750
Number of Degrees of Freedom:  4
R-squared:                    0.7842
Adj R-squared:                0.7833
Rmse:                        0.1141

```

```
F-stat (3, 746): 903.4509, p-value: 0.0000
Degrees of Freedom: model 3, resid 746
```

```
-----Summary of Estimated Coefficients-----
Variable      Coef      Std Err    t-stat    p-value    CI 2.5%    CI 97.5%
-----
return        4.4258     0.3568     12.41     0.0000     3.7265     5.1250
FE_GOOG       -0.1562     0.0102    -15.30     0.0000    -0.1762    -0.1362
FE_MSFT        0.3525     0.0102     34.48     0.0000     0.3325     0.3725
intercept      0.1292     0.0087     14.85     0.0000     0.1121     0.1462
-----End of Summary-----
```

Because we ran the regression with an intercept, one of the dummy variables must be dropped or the design matrix will not be full rank. If we do not use an intercept, all of the dummy variables will be included:

```
In [257]: fe_model = ols(y=volume, x={'return' : np.abs(rets)},
.....:                  entity_effects=True, intercept=False)
.....:
```

```
In [258]: fe_model
```

```
Out[258]:
-----Summary of Regression Analysis-----
Formula: Y ~ <return> + <FE_AAPL> + <FE_GOOG> + <FE_MSFT>
Number of Observations: 750
Number of Degrees of Freedom: 4
R-squared: 0.7842
Adj R-squared: 0.7833
Rmse: 0.1141
F-stat (4, 746): 903.4509, p-value: 0.0000
Degrees of Freedom: model 3, resid 746
-----Summary of Estimated Coefficients-----
Variable      Coef      Std Err    t-stat    p-value    CI 2.5%    CI 97.5%
-----
return        4.4258     0.3568     12.41     0.0000     3.7265     5.1250
FE_AAPL        0.1292     0.0087     14.85     0.0000     0.1121     0.1462
FE_GOOG       -0.0270     0.0085     -3.17     0.0016    -0.0438    -0.0103
FE_MSFT        0.4817     0.0084     57.56     0.0000     0.4653     0.4981
-----End of Summary-----
```

We can also include *time effects*, which demeans the data cross-sectionally at each point in time (equivalent to including dummy variables for each date). More mathematical care must be taken to properly compute the standard errors in this case:

```
In [259]: te_model = ols(y=volume, x={'return' : np.abs(rets)},
.....:                  time_effects=True, entity_effects=True)
.....:
```

```
In [260]: te_model
```

```
Out[260]:
-----Summary of Regression Analysis-----
Formula: Y ~ <return> + <FE_GOOG> + <FE_MSFT>
Number of Observations: 750
Number of Degrees of Freedom: 253
R-squared: 0.8451
Adj R-squared: 0.7666
Rmse: 0.1117
F-stat (3, 497): 10.7603, p-value: 0.0000
Degrees of Freedom: model 252, resid 497
-----Summary of Estimated Coefficients-----
Variable      Coef      Std Err    t-stat    p-value    CI 2.5%    CI 97.5%
-----
```

```
-----  
      return      3.6492      nan      nan      0.0000      nan      nan  
FE_GOOG      -0.1569      nan      nan      0.0000      nan      nan  
FE_MSFT       0.3512      nan      nan      0.0000      nan      nan  
-----End of Summary-----
```

Here the intercept (the mean term) is dropped by default because it will be 0 according to the model assumptions, having subtracted off the group means.

### 8.4.3 Result fields and tests

We'll leave it to the user to explore the docstrings and source, especially as we'll be moving this code into statsmodels in the near future.



# WORKING WITH MISSING DATA

In this section, we will discuss missing (also referred to as NA) values in pandas.

---

**Note:** The choice of using NaN internally to denote missing data was largely for simplicity and performance reasons. It differs from the MaskedArray approach of, for example, `scikits.timeseries`. We are hopeful that NumPy will soon be able to provide a native NA type solution (similar to R) performant enough to be used in pandas.

---

## 9.1 Missing data basics

### 9.1.1 When / why does data become missing?

Some might quibble over our usage of *missing*. By “missing” we simply mean **null** or “not present for whatever reason”. Many data sets simply arrive with missing data, either because it exists and was not collected or it never existed. For example, in a collection of financial time series, some of the time series might start on different dates. Thus, values prior to the start date would generally be marked as missing.

In pandas, one of the most common ways that missing data is **introduced** into a data set is by reindexing. For example

```
In [925]: df = DataFrame(randn(5, 3), index=['a', 'c', 'e', 'f', 'h'],
.....:                  columns=['one', 'two', 'three'])
.....:
```

```
In [926]: df['four'] = 'bar'
```

```
In [927]: df['five'] = df['one'] > 0
```

```
In [928]: df
```

```
Out[928]:
```

	one	two	three	four	five
a	0.059117	1.138469	-2.400634	bar	True
c	-0.280853	0.025653	-1.386071	bar	False
e	0.863937	0.252462	1.500571	bar	True
f	1.053202	-2.338595	-0.374279	bar	True
h	-2.359958	-1.157886	-0.551865	bar	False

```
In [929]: df2 = df.reindex(['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h'])
```

```
In [930]: df2
```

```
Out[930]:
```

	one	two	three	four	five
--	-----	-----	-------	------	------

```
a  0.059117  1.138469 -2.400634  bar  True
b         NaN         NaN         NaN  NaN  NaN
c -0.280853  0.025653 -1.386071  bar  False
d         NaN         NaN         NaN  NaN  NaN
e  0.863937  0.252462  1.500571  bar  True
f  1.053202 -2.338595 -0.374279  bar  True
g         NaN         NaN         NaN  NaN  NaN
h -2.359958 -1.157886 -0.551865  bar  False
```

## 9.1.2 Values considered “missing”

As data comes in many shapes and forms, pandas aims to be flexible with regard to handling missing data. While NaN is the default missing value marker for reasons of computational speed and convenience, we need to be able to easily detect this value with data of different types: floating point, integer, boolean, and general object. In many cases, however, the Python None will arise and we wish to also consider that “missing” or “null”. Lastly, for legacy reasons inf and -inf are also considered to be “null” in computations. Since in NumPy divide-by-zero generates inf or -inf and not NaN, I think you will find this is a worthwhile trade-off (Zen of Python: “practicality beats purity”).

To make detecting missing values easier (and across different array dtypes), pandas provides the isnull() and notnull() functions, which are also methods on Series objects:

```
In [931]: df2['one']
```

```
Out[931]:
a    0.059117
b         NaN
c   -0.280853
d         NaN
e    0.863937
f    1.053202
g         NaN
h   -2.359958
Name: one
```

```
In [932]: isnull(df2['one'])
```

```
Out[932]:
a   False
b    True
c   False
d    True
e   False
f   False
g    True
h   False
Name: one
```

```
In [933]: df2['four'].notnull()
```

```
Out[933]:
a    True
b   False
c    True
d   False
e    True
f    True
g   False
h    True
```

**Summary:** NaN, inf, -inf, and None (in object arrays) are all considered missing by the isnull and notnull functions.

## 9.2 Calculations with missing data

Missing values propagate naturally through arithmetic operations between pandas objects.

```
In [934]: a
```

```
Out[934]:
```

	one	two
a	0.059117	1.138469
b	0.059117	1.138469
c	-0.280853	0.025653
d	-0.280853	0.025653
e	0.863937	0.252462

```
In [935]: b
```

```
Out[935]:
```

	one	two	three
a	0.059117	1.138469	-2.400634
b	NaN	NaN	NaN
c	-0.280853	0.025653	-1.386071
d	NaN	NaN	NaN
e	0.863937	0.252462	1.500571

```
In [936]: a + b
```

```
Out[936]:
```

	one	three	two
a	0.118234	NaN	2.276938
b	NaN	NaN	NaN
c	-0.561707	NaN	0.051306
d	NaN	NaN	NaN
e	1.727874	NaN	0.504923

The descriptive statistics and computational methods discussed in the *data structure overview* (and listed *here* and *here*) are all written to account for missing data. For example:

- When summing data, NA (missing) values will be treated as zero
- If the data are all NA, the result will be NA
- Methods like **cumsum** and **cumprod** ignore NA values, but preserve them in the resulting arrays

```
In [937]: df
```

```
Out[937]:
```

	one	two	three
a	0.059117	1.138469	-2.400634
b	NaN	NaN	NaN
c	-0.280853	0.025653	-1.386071
d	NaN	NaN	NaN
e	0.863937	0.252462	1.500571
f	1.053202	-2.338595	-0.374279
g	NaN	NaN	NaN
h	-2.359958	-1.157886	-0.551865

```
In [938]: df['one'].sum()
```

```
Out[938]: -0.66455558290247652
```

```
In [939]: df.mean(1)
```

```
Out[939]:
```

a	-0.401016
b	NaN

```
c    -0.547090
d         NaN
e     0.872323
f    -0.553224
g         NaN
h    -1.356570
```

```
In [940]: df.cumsum()
```

```
Out[940]:
```

	one	two	three
a	0.059117	1.138469	-2.400634
b	NaN	NaN	NaN
c	-0.221736	1.164122	-3.786705
d	NaN	NaN	NaN
e	0.642200	1.416584	-2.286134
f	1.695403	-0.922011	-2.660413
g	NaN	NaN	NaN
h	-0.664556	-2.079897	-3.212278

### 9.2.1 NA values in GroupBy

NA groups in GroupBy are automatically excluded. This behavior is consistent with R, for example.

## 9.3 Cleaning / filling missing data

pandas objects are equipped with various data manipulation methods for dealing with missing data.

### 9.3.1 Filling missing values: fillna

The `fillna` function can “fill in” NA values with non-null data in a couple of ways, which we illustrate:

#### Replace NA with a scalar value

```
In [941]: df2
```

```
Out[941]:
```

	one	two	three	four	five
a	0.059117	1.138469	-2.400634	bar	True
b	NaN	NaN	NaN	NaN	NaN
c	-0.280853	0.025653	-1.386071	bar	False
d	NaN	NaN	NaN	NaN	NaN
e	0.863937	0.252462	1.500571	bar	True
f	1.053202	-2.338595	-0.374279	bar	True
g	NaN	NaN	NaN	NaN	NaN
h	-2.359958	-1.157886	-0.551865	bar	False

```
In [942]: df2.fillna(0)
```

```
Out[942]:
```

	one	two	three	four	five
a	0.059117	1.138469	-2.400634	bar	True
b	0.000000	0.000000	0.000000	0	0
c	-0.280853	0.025653	-1.386071	bar	False
d	0.000000	0.000000	0.000000	0	0
e	0.863937	0.252462	1.500571	bar	True
f	1.053202	-2.338595	-0.374279	bar	True



```
g 0.000000 0.000000 0.000000 0 0
h -2.359958 -1.157886 -0.551865 bar False
```

```
In [943]: df2['four'].fillna('missing')
```

```
Out[943]:
```

```
a      bar
b  missing
c      bar
d  missing
e      bar
f      bar
g  missing
h      bar
Name: four
```

### Fill gaps forward or backward

Using the same filling arguments as [reindexing](#), we can propagate non-null values forward or backward:

```
In [944]: df
```

```
Out[944]:
```

```
      one      two      three
a  0.059117  1.138469 -2.400634
b      NaN      NaN      NaN
c -0.280853  0.025653 -1.386071
d      NaN      NaN      NaN
e  0.863937  0.252462  1.500571
f  1.053202 -2.338595 -0.374279
g      NaN      NaN      NaN
h -2.359958 -1.157886 -0.551865
```

```
In [945]: df.fillna(method='pad')
```

```
Out[945]:
```

```
      one      two      three
a  0.059117  1.138469 -2.400634
b  0.059117  1.138469 -2.400634
c -0.280853  0.025653 -1.386071
d -0.280853  0.025653 -1.386071
e  0.863937  0.252462  1.500571
f  1.053202 -2.338595 -0.374279
g  1.053202 -2.338595 -0.374279
h -2.359958 -1.157886 -0.551865
```

### Limit the amount of filling

If we only want consecutive gaps filled up to a certain number of data points, we can use the *limit* keyword:

```
In [946]: df
```

```
Out[946]:
```

```
      one      two      three
a  0.059117  1.138469 -2.400634
b      NaN      NaN      NaN
c      NaN      NaN      NaN
d      NaN      NaN      NaN
e  0.863937  0.252462  1.500571
f  1.053202 -2.338595 -0.374279
g      NaN      NaN      NaN
h -2.359958 -1.157886 -0.551865
```

```
In [947]: df.fillna(method='pad', limit=1)
```

```
Out [947]:
```

	one	two	three
a	0.059117	1.138469	-2.400634
b	0.059117	1.138469	-2.400634
c	NaN	NaN	NaN
d	NaN	NaN	NaN
e	0.863937	0.252462	1.500571
f	1.053202	-2.338595	-0.374279
g	1.053202	-2.338595	-0.374279
h	-2.359958	-1.157886	-0.551865

To remind you, these are the available filling methods:

Method	Action
pad / ffill	Fill values forward
bfill / backfill	Fill values backward

With time series data, using pad/ffill is extremely common so that the “last known value” is available at every time point.

### 9.3.2 Dropping axis labels with missing data: dropna

You may wish to simply exclude labels from a data set which refer to missing data. To do this, use the **dropna** method:

```
In [948]: df
```

```
Out [948]:
```

	one	two	three
a	0.059117	1.138469	-2.400634
b	NaN	0.000000	0.000000
c	NaN	0.000000	0.000000
d	NaN	0.000000	0.000000
e	0.863937	0.252462	1.500571
f	1.053202	-2.338595	-0.374279
g	NaN	0.000000	0.000000
h	-2.359958	-1.157886	-0.551865

```
In [949]: df.dropna(axis=0)
```

```
Out [949]:
```

	one	two	three
a	0.059117	1.138469	-2.400634
e	0.863937	0.252462	1.500571
f	1.053202	-2.338595	-0.374279
h	-2.359958	-1.157886	-0.551865

```
In [950]: df.dropna(axis=1)
```

```
Out [950]:
```

	two	three
a	1.138469	-2.400634
b	0.000000	0.000000
c	0.000000	0.000000
d	0.000000	0.000000
e	0.252462	1.500571
f	-2.338595	-0.374279
g	0.000000	0.000000
h	-1.157886	-0.551865

```
In [951]: df['one'].dropna()
```

```
Out [951]:
```

```
a    0.059117
e    0.863937
f    1.053202
h   -2.359958
Name: one
```

**dropna** is presently only implemented for Series and DataFrame, but will be eventually added to Panel. Series.dropna is a simpler method as it only has one axis to consider. DataFrame.dropna has considerably more options, which can be examined *in the API*.

### 9.3.3 Interpolation

A linear **interpolate** method has been implemented on Series. The default interpolation assumes equally spaced points.

```
In [952]: ts.count()
Out[952]: 61
```

```
In [953]: ts.head()
Out[953]:
2000-01-31    0.469112
2000-02-29         NaN
2000-03-31         NaN
2000-04-28         NaN
2000-05-31         NaN
Freq: BM
```

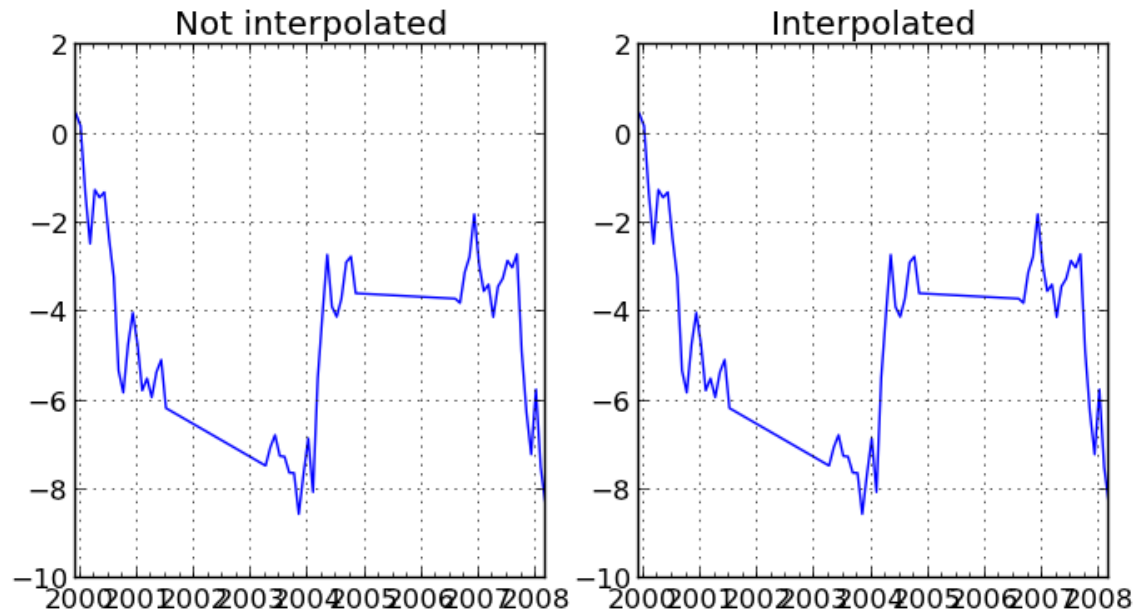
```
In [954]: ts.interpolate().count()
Out[954]: 100
```

```
In [955]: ts.interpolate().head()
Out[955]:
2000-01-31    0.469112
2000-02-29    0.435428
2000-03-31    0.401743
2000-04-28    0.368059
2000-05-31    0.334374
Freq: BM
```

```
In [956]: fig = plt.figure()
```

```
In [957]: ts.interpolate().plot()
Out[957]: <matplotlib.axes.AxesSubplot at 0xbc9a790>
```

```
In [958]: plt.close('all')
```



Index aware interpolation is available via the `method` keyword:

```
In [959]: ts
Out[959]:
2000-01-31    0.469112
2000-02-29         NaN
2002-07-31   -5.689738
2005-01-31         NaN
2008-04-30   -8.916232
```

```
In [960]: ts.interpolate()
Out[960]:
2000-01-31    0.469112
2000-02-29   -2.610313
2002-07-31   -5.689738
2005-01-31   -7.302985
2008-04-30   -8.916232
```

```
In [961]: ts.interpolate(method='time')
Out[961]:
2000-01-31    0.469112
2000-02-29    0.273272
2002-07-31   -5.689738
2005-01-31   -7.095568
2008-04-30   -8.916232
```

For a floating-point index, use `method='values'`:

```
In [962]: ser
Out[962]:
0      0
1     NaN
10     10

In [963]: ser.interpolate()
Out[963]:
0      0
```

```
1      5
10     10
```

```
In [964]: ser.interpolate(method='values')
```

```
Out[964]:
```

```
0      0
1      1
10     10
```

### 9.3.4 Replacing Generic Values

Often times we want to replace arbitrary values with other values. New in v0.8 is the `replace` method in `Series/DataFrame` that provides an efficient yet flexible way to perform such replacements.

For a `Series`, you can replace a single value or a list of values by another value:

```
In [965]: ser = Series([0., 1., 2., 3., 4.])
```

```
In [966]: ser.replace(0, 5)
```

```
Out[966]:
```

```
0      5
1      1
2      2
3      3
4      4
```

You can replace a list of values by a list of other values:

```
In [967]: ser.replace([0, 1, 2, 3, 4], [4, 3, 2, 1, 0])
```

```
Out[967]:
```

```
0      4
1      3
2      2
3      1
4      0
```

You can also specify a mapping dict:

```
In [968]: ser.replace({0: 10, 1: 100})
```

```
Out[968]:
```

```
0      10
1     100
2      2
3      3
4      4
```

For a `DataFrame`, you can specify individual values by column:

```
In [969]: df = DataFrame({'a': [0, 1, 2, 3, 4], 'b': [5, 6, 7, 8, 9]})
```

```
In [970]: df.replace({'a': 0, 'b': 5}, 100)
```

```
Out[970]:
```

```
   a  b
0 100 100
1   1   6
2   2   7
3   3   8
4   4   9
```

Instead of replacing with specified values, you can treat all given values as missing and interpolate over them:

```
In [971]: ser.replace([1, 2, 3], method='pad')
Out[971]:
0    0
1    0
2    0
3    0
4    4
```

## 9.4 Missing data casting rules and indexing

While pandas supports storing arrays of integer and boolean type, these types are not capable of storing missing data. Until we can switch to using a native NA type in NumPy, we've established some “casting rules” when reindexing will cause missing data to be introduced into, say, a Series or DataFrame. Here they are:

data type	Cast to
integer	float
boolean	object
float	no cast
object	no cast

For example:

```
In [972]: s = Series(randn(5), index=[0, 2, 4, 6, 7])
```

```
In [973]: s > 0
```

```
Out[973]:
0    False
2     True
4     True
6     True
7     True
```

```
In [974]: (s > 0).dtype
```

```
Out[974]: dtype('bool')
```

```
In [975]: crit = (s > 0).reindex(range(8))
```

```
In [976]: crit
```

```
Out[976]:
0    False
1     NaN
2     True
3     NaN
4     True
5     NaN
6     True
7     True
```

```
In [977]: crit.dtype
```

```
Out[977]: dtype('object')
```

Ordinarily NumPy will complain if you try to use an object array (even if it contains boolean values) instead of a boolean array to get or set values from an ndarray (e.g. selecting values based on some criteria). If a boolean vector contains NAs, an exception will be generated:

```
In [978]: reindexed = s.reindex(range(8)).fillna(0)
```

```
In [979]: reindexed[crit]
```

```
-----
ValueError                                Traceback (most recent call last)
<ipython-input-979-2da204ed1ac7> in <module>()
----> 1 reindexed[crit]
/home/wesm/code/pandas/pandas/core/series.pyc in __getitem__(self, key)
    447         # special handling of boolean data with NAs stored in object
    448         # arrays. Since we can't represent NA with dtype=bool
--> 449         if _is_bool_indexer(key):
    450             key = self._check_bool_indexer(key)
    451             key = np.asarray(key, dtype=bool)
/home/wesm/code/pandas/pandas/core/common.pyc in _is_bool_indexer(key)
    499         if not lib.is_bool_array(key):
    500             if isnull(key).any():
--> 501                 raise ValueError('cannot index with vector containing '
    502                                 'NA / NaN values')
    503         return False
ValueError: cannot index with vector containing NA / NaN values
```

However, these can be filled in using **fillna** and it will work fine:

```
In [980]: reindexed[crit.fillna(False)]
```

```
Out[980]:
2    1.314232
4    0.690579
6    0.995761
7    2.396780
```

```
In [981]: reindexed[crit.fillna(True)]
```

```
Out[981]:
1    0.000000
2    1.314232
3    0.000000
4    0.690579
5    0.000000
6    0.995761
7    2.396780
```





# GROUP BY: SPLIT-APPLY-COMBINE

By “group by” we refer to a process involving one or more of the following steps

- **Splitting** the data into groups based on some criteria
- **Applying** a function to each group independently
- **Combining** the results into a data structure

Of these, the split step is the most straightforward. In fact, in many situations you may wish to split the data set into groups and do something with those groups yourself. In the apply step, we might wish to one of the following:

- **Aggregation:** computing a summary statistic (or statistics) about each group. Some examples:
  - Compute group sums or means
  - Compute group sizes / counts
- **Transformation:** perform some group-specific computations and return a like-indexed. Some examples:
  - Standardizing data (zscore) within group
  - Filling NAs within groups with a value derived from each group
- Some combination of the above: GroupBy will examine the results of the apply step and try to return a sensibly combined result if it doesn't fit into either of the above two categories

Since the set of object instance method on pandas data structures are generally rich and expressive, we often simply want to invoke, say, a DataFrame function on each group. The name GroupBy should be quite familiar to those who have used a SQL-based tool (or `itertools`), in which you can write code like:

```
SELECT Column1, Column2, mean(Column3), sum(Column4)
FROM SomeTable
GROUP BY Column1, Column2
```

We aim to make operations like this natural and easy to express using pandas. We'll address each area of GroupBy functionality then provide some non-trivial examples / use cases.

## 10.1 Splitting an object into groups

pandas objects can be split on any of their axes. The abstract definition of grouping is to provide a mapping of labels to group names. To create a GroupBy object (more on what the GroupBy object is later), you do the following:

```
>>> grouped = obj.groupby(key)
>>> grouped = obj.groupby(key, axis=1)
>>> grouped = obj.groupby([key1, key2])
```

The mapping can be specified many different ways:

- A Python function, to be called on each of the axis labels
- A list or NumPy array of the same length as the selected axis
- A dict or Series, providing a label → group name mapping
- For DataFrame objects, a string indicating a column to be used to group. Of course `df.groupby('A')` is just syntactic sugar for `df.groupby(df['A'])`, but it makes life simpler
- A list of any of the above things

Collectively we refer to the grouping objects as the **keys**. For example, consider the following DataFrame:

```
In [444]: df = DataFrame({'A' : ['foo', 'bar', 'foo', 'bar',  
.....:                        'foo', 'bar', 'foo', 'foo'],  
.....:                  'B' : ['one', 'one', 'two', 'three',  
.....:                        'two', 'two', 'one', 'three'],  
.....:                  'C' : randn(8), 'D' : randn(8)})  
.....:
```

```
In [445]: df
```

```
Out[445]:
```

	A	B	C	D
0	foo	one	0.469112	-0.861849
1	bar	one	-0.282863	-2.104569
2	foo	two	-1.509059	-0.494929
3	bar	three	-1.135632	1.071804
4	foo	two	1.212112	0.721555
5	bar	two	-0.173215	-0.706771
6	foo	one	0.119209	-1.039575
7	foo	three	-1.044236	0.271860

We could naturally group by either the A or B columns or both:

```
In [446]: grouped = df.groupby('A')
```

```
In [447]: grouped = df.groupby(['A', 'B'])
```

These will split the DataFrame on its index (rows). We could also split by the columns:

```
In [448]: def get_letter_type(letter):  
.....:     if letter.lower() in 'aeiou':  
.....:         return 'vowel'  
.....:     else:  
.....:         return 'consonant'  
.....:
```

```
In [449]: grouped = df.groupby(get_letter_type, axis=1)
```

Starting with 0.8, pandas Index objects now supports duplicate values. If a non-unique index is used as the group key in a groupby operation, all values for the same index value will be considered to be in one group and thus the output of aggregation functions will only contain unique index values:

```
In [450]: lst = [1, 2, 3, 1, 2, 3]
```

```
In [451]: s = Series([1, 2, 3, 10, 20, 30], lst)
```

```
In [452]: grouped = s.groupby(level=0)
```

```
In [453]: grouped.first()
```

```
Out [453]:
```

```
1    1
2    2
3    3
```

```
In [454]: grouped.last()
```

```
Out [454]:
```

```
1    10
2    20
3    30
```

```
In [455]: grouped.sum()
```

```
Out [455]:
```

```
1    11
2    22
3    33
```

Note that **no splitting occurs** until it's needed. Creating the GroupBy object only verifies that you've passed a valid mapping.

---

**Note:** Many kinds of complicated data manipulations can be expressed in terms of GroupBy operations (though can't be guaranteed to be the most efficient). You can get quite creative with the label mapping functions.

---

### 10.1.1 GroupBy object attributes

The `groups` attribute is a dict whose keys are the computed unique groups and corresponding values being the axis labels belonging to each group. In the above example we have:

```
In [456]: df.groupby('A').groups
```

```
Out [456]: {'bar': [1, 3, 5], 'foo': [0, 2, 4, 6, 7]}
```

```
In [457]: df.groupby(get_letter_type, axis=1).groups
```

```
Out [457]: {'consonant': ['B', 'C', 'D'], 'vowel': ['A']}
```

Calling the standard Python `len` function on the GroupBy object just returns the length of the groups dict, so it is largely just a convenience:

```
In [458]: grouped = df.groupby(['A', 'B'])
```

```
In [459]: grouped.groups
```

```
Out [459]:
```

```
{('bar', 'one'): [1],
 ('bar', 'three'): [3],
 ('bar', 'two'): [5],
 ('foo', 'one'): [0, 6],
 ('foo', 'three'): [7],
 ('foo', 'two'): [2, 4]}
```

```
In [460]: len(grouped)
```

```
Out [460]: 6
```

By default the group keys are sorted during the groupby operation. You may however pass `sort=False` for potential speedups:

```
In [461]: df2 = DataFrame({'X' : ['B', 'B', 'A', 'A'], 'Y' : [1, 2, 3, 4]})
```

```
In [462]: df2.groupby(['X'], sort=True).sum()
```

```
Out[462]:
```

```
      Y
X
A    7
B    3
```

```
In [463]: df2.groupby(['X'], sort=False).sum()
```

```
Out[463]:
```

```
      Y
X
B    3
A    7
```

### 10.1.2 GroupBy with MultiIndex

With *hierarchically-indexed data*, it's quite natural to group by one of the levels of the hierarchy.

```
In [464]: s
```

```
Out[464]:
```

```
first second
bar    one    -0.424972
      two     0.567020
baz    one     0.276232
      two    -1.087401
foo    one    -0.673690
      two     0.113648
qux    one    -1.478427
      two     0.524988
```

```
In [465]: grouped = s.groupby(level=0)
```

```
In [466]: grouped.sum()
```

```
Out[466]:
```

```
first
bar    0.142048
baz   -0.811169
foo   -0.560041
qux   -0.953439
```

If the MultiIndex has names specified, these can be passed instead of the level number:

```
In [467]: s.groupby(level='second').sum()
```

```
Out[467]:
```

```
second
one    -2.300857
two     0.118256
```

The aggregation functions such as `sum` will take the `level` parameter directly. Additionally, the resulting index will be named according to the chosen level:

```
In [468]: s.sum(level='second')
```

```
Out[468]:
```

```
second
one    -2.300857
two     0.118256
```

Also as of v0.6, grouping with multiple levels is supported.

```
In [469]: s
Out[469]:
first  second  third
bar    doo     one    0.404705
        two    0.577046
baz    bee     one   -1.715002
        two   -1.039268
foo    bop     one   -0.370647
        two   -1.157892
qux    bop     one   -1.344312
        two    0.844885

In [470]: s.groupby(level=['first', 'second']).sum()
Out[470]:
first  second
bar    doo     0.981751
baz    bee    -2.754270
foo    bop    -1.528539
qux    bop    -0.499427
```

More on the `sum` function and aggregation later.

### 10.1.3 DataFrame column selection in GroupBy

Once you have created the `GroupBy` object from a `DataFrame`, for example, you might want to do something different for each of the columns. Thus, using `[]` similar to getting a column from a `DataFrame`, you can do:

```
In [471]: grouped = df.groupby(['A'])

In [472]: grouped_C = grouped['C']

In [473]: grouped_D = grouped['D']
```

This is mainly syntactic sugar for the alternative and much more verbose:

```
In [474]: df['C'].groupby(df['A'])
Out[474]: <pandas.core.groupby.SeriesGroupBy at 0x9289090>
```

Additionally this method avoids recomputing the internal grouping information derived from the passed key.

## 10.2 Iterating through groups

With the `GroupBy` object in hand, iterating through the grouped data is very natural and functions similarly to `itertools.groupby`:

```
In [475]: grouped = df.groupby('A')

In [476]: for name, group in grouped:
.....:     print name
.....:     print group
.....:
bar
   A      B      C      D
1 bar  one -0.282863 -2.104569
3 bar three -1.135632  1.071804
```

```
5 bar    two -0.173215 -0.706771
foo
   A      B      C      D
0 foo    one  0.469112 -0.861849
2 foo    two -1.509059 -0.494929
4 foo    two  1.212112  0.721555
6 foo    one  0.119209 -1.039575
7 foo   three -1.044236  0.271860
```

In the case of grouping by multiple keys, the group name will be a tuple:

```
In [477]: for name, group in df.groupby(['A', 'B']):
.....:     print name
.....:     print group
.....:
('bar', 'one')
   A      B      C      D
1 bar    one -0.282863 -2.104569
('bar', 'three')
   A      B      C      D
3 bar   three -1.135632  1.071804
('bar', 'two')
   A      B      C      D
5 bar    two -0.173215 -0.706771
('foo', 'one')
   A      B      C      D
0 foo    one  0.469112 -0.861849
6 foo    one  0.119209 -1.039575
('foo', 'three')
   A      B      C      D
7 foo   three -1.044236  0.27186
('foo', 'two')
   A      B      C      D
2 foo    two -1.509059 -0.494929
4 foo    two  1.212112  0.721555
```

It's standard Python-fu but remember you can unpack the tuple in the for loop statement if you wish: `for (k1, k2), group in grouped:`.

## 10.3 Aggregation

Once the `GroupBy` object has been created, several methods are available to perform a computation on the grouped data. An obvious one is aggregation via the `aggregate` or equivalently `agg` method:

```
In [478]: grouped = df.groupby('A')
```

```
In [479]: grouped.agg(np.sum)
```

```
Out[479]:
           C      D
A
bar -1.591710 -1.739537
foo -0.752861 -1.402938
```

```
In [480]: grouped = df.groupby(['A', 'B'])
```

```
In [481]: grouped.agg(np.sum)
```

```
Out[481]:
```

		C	D
A	B		
bar	one	-0.282863	-2.104569
	three	-1.135632	1.071804
	two	-0.173215	-0.706771
foo	one	0.588321	-1.901424
	three	-1.044236	0.271860
	two	-0.296946	0.226626

As you can see, the result of the aggregation will have the group names as the new index along the grouped axis. In the case of multiple keys, the result is a *MultiIndex* by default, though this can be changed by using the `as_index` option:

```
In [482]: grouped = df.groupby(['A', 'B'], as_index=False)
```

```
In [483]: grouped.agg(np.sum)
```

```
Out[483]:
```

	A	B	C	D
0	bar	one	-0.282863	-2.104569
1	bar	three	-1.135632	1.071804
2	bar	two	-0.173215	-0.706771
3	foo	one	0.588321	-1.901424
4	foo	three	-1.044236	0.271860
5	foo	two	-0.296946	0.226626

```
In [484]: df.groupby('A', as_index=False).sum()
```

```
Out[484]:
```

	A	C	D
0	bar	-1.591710	-1.739537
1	foo	-0.752861	-1.402938

Note that you could use the `reset_index` *DataFrame* function to achieve the same result as the column names are stored in the resulting *MultiIndex*:

```
In [485]: df.groupby(['A', 'B']).sum().reset_index()
```

```
Out[485]:
```

	A	B	C	D
0	bar	one	-0.282863	-2.104569
1	bar	three	-1.135632	1.071804
2	bar	two	-0.173215	-0.706771
3	foo	one	0.588321	-1.901424
4	foo	three	-1.044236	0.271860
5	foo	two	-0.296946	0.226626

Another simple aggregation example is to compute the size of each group. This is included in *GroupBy* as the `size` method. It returns a *Series* whose index are the group names and whose values are the sizes of each group.

```
In [486]: grouped.size()
```

```
Out[486]:
```

A	B	
bar	one	1
	three	1
	two	1
foo	one	2
	three	1
	two	2

### 10.3.1 Applying multiple functions at once

With grouped Series you can also pass a list or dict of functions to do aggregation with, outputting a DataFrame:

```
In [487]: grouped = df.groupby('A')

In [488]: grouped['C'].agg([np.sum, np.mean, np.std])
Out[488]:
```

	sum	mean	std
A			
bar	-1.591710	-0.530570	0.526860
foo	-0.752861	-0.150572	1.113308

If a dict is passed, the keys will be used to name the columns. Otherwise the function's name (stored in the function object) will be used.

```
In [489]: grouped['D'].agg({'result1' : np.sum,
.....:                    'result2' : np.mean})
.....:
Out[489]:
```

	result2	result1
A		
bar	-0.579846	-1.739537
foo	-0.280588	-1.402938

On a grouped DataFrame, you can pass a list of functions to apply to each column, which produces an aggregated result with a hierarchical index:

```
In [490]: grouped.agg([np.sum, np.mean, np.std])
Out[490]:
```

	C			D		
	sum	mean	std	sum	mean	std
A						
bar	-1.591710	-0.530570	0.526860	-1.739537	-0.579846	1.591986
foo	-0.752861	-0.150572	1.113308	-1.402938	-0.280588	0.753219

Passing a dict of functions has different behavior by default, see the next section.

### 10.3.2 Applying different functions to DataFrame columns

By passing a dict to `aggregate` you can apply a different aggregation to the columns of a DataFrame:

```
In [491]: grouped.agg({'C' : np.sum,
.....:                'D' : lambda x: np.std(x, ddof=1)})
.....:
Out[491]:
```

	C	D
A		
bar	-1.591710	1.591986
foo	-0.752861	0.753219

The function names can also be strings. In order for a string to be valid it must be either implemented on `GroupBy` or available via *dispatching*:

```
In [492]: grouped.agg({'C' : 'sum', 'D' : 'std'})
Out[492]:
```

	C	D
A		



```
bar -1.591710  1.591986
foo -0.752861  0.753219
```

### 10.3.3 Cython-optimized aggregation functions

Some common aggregations, currently only `sum`, `mean`, and `std`, have optimized Cython implementations:

```
In [493]: df.groupby('A').sum()
Out[493]:
```

```
      C      D
A
bar -1.591710 -1.739537
foo -0.752861 -1.402938
```

```
In [494]: df.groupby(['A', 'B']).mean()
Out[494]:
```

```
      C      D
A  B
bar one  -0.282863 -2.104569
     three -1.135632  1.071804
     two  -0.173215 -0.706771
foo one   0.294161 -0.950712
     three -1.044236  0.271860
     two  -0.148473  0.113313
```

Of course `sum` and `mean` are implemented on pandas objects, so the above code would work even without the special versions via dispatching (see below).

## 10.4 Transformation

The `transform` method returns an object that is indexed the same (same size) as the one being grouped. Thus, the passed transform function should return a result that is the same size as the group chunk. For example, suppose we wished to standardize the data within each group:

```
In [495]: index = date_range('10/1/1999', periods=1100)
```

```
In [496]: ts = Series(np.random.normal(0.5, 2, 1100), index)
```

```
In [497]: ts = rolling_mean(ts, 100, 100).dropna()
```

```
In [498]: ts.head()
```

```
Out[498]:
2000-01-08    0.536925
2000-01-09    0.494448
2000-01-10    0.496114
2000-01-11    0.443475
2000-01-12    0.474744
Freq: D
```

```
In [499]: ts.tail()
```

```
Out[499]:
2002-09-30    0.978859
2002-10-01    0.994704
2002-10-02    0.953789
2002-10-03    0.932345
```

```
2002-10-04    0.915581
Freq: D
```

```
In [500]: key = lambda x: x.year
```

```
In [501]: zscore = lambda x: (x - x.mean()) / x.std()
```

```
In [502]: transformed = ts.groupby(key).transform(zscore)
```

We would expect the result to now have mean 0 and standard deviation 1 within each group, which we can easily check:

```
# Original Data
```

```
In [503]: grouped = ts.groupby(key)
```

```
In [504]: grouped.mean()
```

```
Out[504]:
2000    0.416344
2001    0.416987
2002    0.599380
```

```
In [505]: grouped.std()
```

```
Out[505]:
2000    0.174755
2001    0.309640
2002    0.266172
```

```
# Transformed Data
```

```
In [506]: grouped_trans = transformed.groupby(key)
```

```
In [507]: grouped_trans.mean()
```

```
Out[507]:
2000    -0
2001    -0
2002    -0
```

```
In [508]: grouped_trans.std()
```

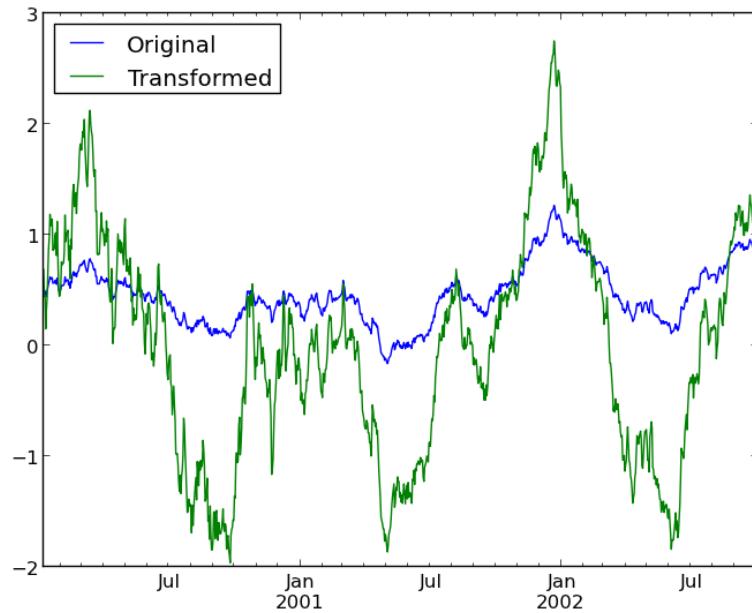
```
Out[508]:
2000    1
2001    1
2002    1
```

We can also visually compare the original and transformed data sets.

```
In [509]: compare = DataFrame({'Original': ts, 'Transformed': transformed})
```

```
In [510]: compare.plot()
```

```
Out[510]: <matplotlib.axes.AxesSubplot at 0xb07fd90>
```



Another common data transform is to replace missing data with the group mean.

```
In [511]: data_df
Out[511]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1000 entries, 0 to 999
Data columns:
A      908 non-null values
B      953 non-null values
C      820 non-null values
dtypes: float64(3)

In [512]: countries = np.array(['US', 'UK', 'GR', 'JP'])

In [513]: key = countries[np.random.randint(0, 4, 1000)]

In [514]: grouped = data_df.groupby(key)

# Non-NA count in each group
In [515]: grouped.count()
Out[515]:
      A      B      C
GR  219  223  194
JP  238  250  211
UK  228  239  213
US  223  241  202
```

```
In [516]: f = lambda x: x.fillna(x.mean())

In [517]: transformed = grouped.transform(f)
```

We can verify that the group means have not changed in the transformed data and that the transformed data contains no NAs.

```
In [518]: grouped_trans = transformed.groupby(key)

In [519]: grouped.mean() # original group means
```

```
Out[519]:
```

	A	B	C
GR	0.093655	-0.004978	-0.049883
JP	-0.067605	0.025828	0.006752
UK	-0.054246	0.031742	0.068974
US	0.084334	-0.013433	0.056589

```
In [520]: grouped_trans.mean() # transformation did not change group means
```

```
Out[520]:
```

	A	B	C
GR	0.093655	-0.004978	-0.049883
JP	-0.067605	0.025828	0.006752
UK	-0.054246	0.031742	0.068974
US	0.084334	-0.013433	0.056589

```
In [521]: grouped.count() # original has some missing data points
```

```
Out[521]:
```

	A	B	C
GR	219	223	194
JP	238	250	211
UK	228	239	213
US	223	241	202

```
In [522]: grouped_trans.count() # counts after transformation
```

```
Out[522]:
```

	A	B	C
GR	234	234	234
JP	264	264	264
UK	251	251	251
US	251	251	251

```
In [523]: grouped_trans.size() # Verify non-NA count equals group size
```

```
Out[523]:
```

GR	234
JP	264
UK	251
US	251

## 10.5 Dispatching to instance methods

When doing an aggregation or transformation, you might just want to call an instance method on each data group. This is pretty easy to do by passing lambda functions:

```
In [524]: grouped = df.groupby('A')
```

```
In [525]: grouped.agg(lambda x: x.std())
```

```
Out[525]:
```

	B	C	D
A			
bar	NaN	0.526860	1.591986
foo	NaN	1.113308	0.753219

But, it's rather verbose and can be untidy if you need to pass additional arguments. Using a bit of metaprogramming cleverness, GroupBy now has the ability to “dispatch” method calls to the groups:

```
In [526]: grouped.std()
```

```
Out[526]:
```

	C	D
A		
bar	0.526860	1.591986
foo	1.113308	0.753219

What is actually happening here is that a function wrapper is being generated. When invoked, it takes any passed arguments and invokes the function with any arguments on each group (in the above example, the `std` function). The results are then combined together much in the style of `agg` and `transform` (it actually uses `apply` to infer the grouping, documented next). This enables some operations to be carried out rather succinctly:

```
In [527]: tsdf = DataFrame(randn(1000, 3),
.....:                    index=date_range('1/1/2000', periods=1000),
.....:                    columns=['A', 'B', 'C'])
.....:
```

```
In [528]: tsdf.ix[:,2] = np.nan
```

```
In [529]: grouped = tsdf.groupby(lambda x: x.year)
```

```
In [530]: grouped.fillna(method='pad')
```

```
Out[530]:
```

```
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 1000 entries, 2000-01-01 00:00:00 to 2002-09-26 00:00:00
Freq: D
Data columns:
A      998 non-null values
B      998 non-null values
C      998 non-null values
dtypes: float64(3)
```

In this example, we chopped the collection of time series into yearly chunks then independently called *fillna* on the groups.

## 10.6 Flexible apply

Some operations on the grouped data might not fit into either the aggregate or transform categories. Or, you may simply want `GroupBy` to infer how to combine the results. For these, use the `apply` function, which can be substituted for both `aggregate` and `transform` in many standard use cases. However, `apply` can handle some exceptional use cases, for example:

```
In [531]: df
```

```
Out[531]:
```

	A	B	C	D
0	foo	one	0.469112	-0.861849
1	bar	one	-0.282863	-2.104569
2	foo	two	-1.509059	-0.494929
3	bar	three	-1.135632	1.071804
4	foo	two	1.212112	0.721555
5	bar	two	-0.173215	-0.706771
6	foo	one	0.119209	-1.039575
7	foo	three	-1.044236	0.271860

```
In [532]: grouped = df.groupby('A')
```

```
# could also just call .describe()
In [533]: grouped['C'].apply(lambda x: x.describe())
Out[533]:
A
bar  count    3.000000
     mean   -0.530570
     std    0.526860
     min   -1.135632
     25%   -0.709248
     50%   -0.282863
     75%   -0.228039
     max   -0.173215
foo  count    5.000000
     mean   -0.150572
     std    1.113308
     min   -1.509059
     25%   -1.044236
     50%    0.119209
     75%    0.469112
     max    1.212112
```

The dimension of the returned result can also change:

```
In [534]: grouped = df.groupby('A')['C']

In [535]: def f(group):
.....:     return DataFrame({'original' : group,
.....:                      'demeaned' : group - group.mean()})
.....:

In [536]: grouped.apply(f)
Out[536]:
   demeaned  original
0  0.619685  0.469112
1  0.247707 -0.282863
2 -1.358486 -1.509059
3 -0.605062 -1.135632
4  1.362684  1.212112
5  0.357355 -0.173215
6  0.269781  0.119209
7 -0.893664 -1.044236
```

## 10.7 Other useful features

### 10.7.1 Automatic exclusion of “nuisance” columns

Again consider the example DataFrame we’ve been looking at:

```
In [537]: df
Out[537]:
   A      B      C      D
0  foo  one  0.469112 -0.861849
1  bar  one -0.282863 -2.104569
2  foo  two -1.509059 -0.494929
3  bar three -1.135632  1.071804
4  foo  two  1.212112  0.721555
```

```
5 bar    two -0.173215 -0.706771
6 foo    one  0.119209 -1.039575
7 foo  three -1.044236  0.271860
```

Supposed we wished to compute the standard deviation grouped by the A column. There is a slight problem, namely that we don't care about the data in column B. We refer to this as a “nuisance” column. If the passed aggregation function can't be applied to some columns, the troublesome columns will be (silently) dropped. Thus, this does not pose any problems:

```
In [538]: df.groupby('A').std()
```

```
Out [538]:
```

	C	D
A		
bar	0.526860	1.591986
foo	1.113308	0.753219

## 10.7.2 NA group handling

If there are any NaN values in the grouping key, these will be automatically excluded. So there will never be an “NA group”. This was not the case in older versions of pandas, but users were generally discarding the NA group anyway (and supporting it was an implementation headache).

## 10.7.3 Grouping with ordered factors

Categorical variables represented as instance of pandas's `Factor` class can be used as group keys. If so, the order of the levels will be preserved:

```
In [539]: data = Series(np.random.randn(100))
```

```
In [540]: factor = qcut(data, [0, .25, .5, .75, 1.])
```

```
In [541]: data.groupby(factor).mean()
```

```
Out [541]:
```

[-3.713, -0.815]	-1.432664
(-0.815, 0.105]	-0.306693
(0.105, 0.609]	0.356789
(0.609, 2.154]	1.314491





# MERGE, JOIN, AND CONCATENATE

pandas provides various facilities for easily combining together Series, DataFrame, and Panel objects with various kinds of set logic for the indexes and relational algebra functionality in the case of join / merge-type operations.

## 11.1 Concatenating objects

The `concat` function (in the main pandas namespace) does all of the heavy lifting of performing concatenation operations along an axis while performing optional set logic (union or intersection) of the indexes (if any) on the other axes. Note that I say “if any” because there is only a single possible axis of concatenation for Series.

Before diving into all of the details of `concat` and what it can do, here is a simple example:

```
In [821]: df = DataFrame(np.random.randn(10, 4))
```

```
In [822]: df
```

```
Out[822]:
```

	0	1	2	3
0	0.469112	-0.282863	-1.509059	-1.135632
1	1.212112	-0.173215	0.119209	-1.044236
2	-0.861849	-2.104569	-0.494929	1.071804
3	0.721555	-0.706771	-1.039575	0.271860
4	-0.424972	0.567020	0.276232	-1.087401
5	-0.673690	0.113648	-1.478427	0.524988
6	0.404705	0.577046	-1.715002	-1.039268
7	-0.370647	-1.157892	-1.344312	0.844885
8	1.075770	-0.109050	1.643563	-1.469388
9	0.357021	-0.674600	-1.776904	-0.968914

```
# break it into pieces
```

```
In [823]: pieces = [df[:3], df[3:7], df[7:]]
```

```
In [824]: concatenated = concat(pieces)
```

```
In [825]: concatenated
```

```
Out[825]:
```

	0	1	2	3
0	0.469112	-0.282863	-1.509059	-1.135632
1	1.212112	-0.173215	0.119209	-1.044236
2	-0.861849	-2.104569	-0.494929	1.071804
3	0.721555	-0.706771	-1.039575	0.271860
4	-0.424972	0.567020	0.276232	-1.087401
5	-0.673690	0.113648	-1.478427	0.524988
6	0.404705	0.577046	-1.715002	-1.039268

```
7 -0.370647 -1.157892 -1.344312  0.844885
8  1.075770 -0.109050  1.643563 -1.469388
9  0.357021 -0.674600 -1.776904 -0.968914
```

Like its sibling function on ndarrays, `numpy.concatenate`, `pandas.concat` takes a list or dict of homogeneously-typed objects and concatenates them with some configurable handling of “what to do with the other axes”:

```
concat(objs, axis=0, join='outer', join_axes=None, ignore_index=False,
       keys=None, levels=None, names=None, verify_integrity=False)
```

- `objs`: list or dict of Series, DataFrame, or Panel objects. If a dict is passed, the sorted keys will be used as the `keys` argument, unless it is passed, in which case the values will be selected (see below)
- `axis`: {0, 1, ...}, default 0. The axis to concatenate along
- `join`: {'inner', 'outer'}, default 'outer'. How to handle indexes on other axis(es). Outer for union and inner for intersection
- `join_axes`: list of Index objects. Specific indexes to use for the other `n - 1` axes instead of performing inner/outer set logic
- `keys`: sequence, default None. Construct hierarchical index using the passed keys as the outermost level. If multiple levels passed, should contain tuples.
- `levels`: list of sequences, default None. If keys passed, specific levels to use for the resulting MultiIndex. Otherwise they will be inferred from the keys
- `names`: list, default None. Names for the levels in the resulting hierarchical index
- `verify_integrity`: boolean, default False. Check whether the new concatenated axis contains duplicates. This can be very expensive relative to the actual data concatenation
- `ignore_index`: boolean, default False. If True, do not use the index values on the concatenation axis. The resulting axis will be labeled 0, ..., `n - 1`. This is useful if you are concatenating objects where the concatenation axis does not have meaningful indexing information.

Without a little bit of context and example many of these arguments don't make much sense. Let's take the above example. Suppose we wanted to associate specific keys with each of the pieces of the chopped up DataFrame. We can do this using the `keys` argument:

```
In [826]: concatenated = concat(pieces, keys=['first', 'second', 'third'])
```

```
In [827]: concatenated
```

```
Out[827]:
```

```
          0          1          2          3
first 0  0.469112 -0.282863 -1.509059 -1.135632
      1  1.212112 -0.173215  0.119209 -1.044236
      2 -0.861849 -2.104569 -0.494929  1.071804
second 3  0.721555 -0.706771 -1.039575  0.271860
      4 -0.424972  0.567020  0.276232 -1.087401
      5 -0.673690  0.113648 -1.478427  0.524988
      6  0.404705  0.577046 -1.715002 -1.039268
third  7 -0.370647 -1.157892 -1.344312  0.844885
      8  1.075770 -0.109050  1.643563 -1.469388
      9  0.357021 -0.674600 -1.776904 -0.968914
```

As you can see (if you've read the rest of the documentation), the resulting object's index has a *hierarchical index*. This means that we can now do stuff like select out each chunk by key:

```
In [828]: concatenated.ix['second']
Out[828]:
```

	0	1	2	3
3	0.721555	-0.706771	-1.039575	0.271860
4	-0.424972	0.567020	0.276232	-1.087401
5	-0.673690	0.113648	-1.478427	0.524988
6	0.404705	0.577046	-1.715002	-1.039268

It's not a stretch to see how this can be very useful. More detail on this functionality below.

### 11.1.1 Set logic on the other axes

When gluing together multiple DataFrames (or Panels or...), for example, you have a choice of how to handle the other axes (other than the one being concatenated). This can be done in three ways:

- Take the (sorted) union of them all, `join='outer'`. This is the default option as it results in zero information loss.
- Take the intersection, `join='inner'`.
- Use a specific index (in the case of DataFrame) or indexes (in the case of Panel or future higher dimensional objects), i.e. the `join_axes` argument

Here is a example of each of these methods. First, the default `join='outer'` behavior:

```
In [829]: from pandas.util.testing import randn
```

```
In [830]: df = DataFrame(np.random.randn(10, 4), columns=['a', 'b', 'c', 'd'],
.....:                  index=[randn(5) for _ in xrange(10)])
.....:
```

```
In [831]: df
```

```
Out[831]:
```

	a	b	c	d
2dY1o	-1.294524	0.413738	0.276662	-0.472035
eW4ge	-0.013960	-0.362543	-0.006154	-0.923061
fz64l	0.895717	0.805244	-1.206412	2.565646
a60mV	1.431256	1.340309	-1.170299	-0.226169
0tnT9	0.410835	0.813850	0.132003	-0.827317
nJL3w	-0.076467	-1.187678	1.130127	-1.436737
fzroS	-1.413681	1.607920	1.024180	0.569605
5pSvW	0.875906	-2.211372	0.974466	-2.006747
xCmnX	-0.410001	-0.078638	0.545952	-1.219217
MGwb4	-1.226825	0.769804	-1.281247	-0.727707

```
In [832]: concat([df.ix[:7, ['a', 'b']], df.ix[2:-2, ['c']],
.....:            df.ix[-7:, ['d']], axis=1)
.....:
```

```
Out[832]:
```

	a	b	c	d
0tnT9	0.410835	0.813850	0.132003	-0.827317
2dY1o	-1.294524	0.413738	NaN	NaN
5pSvW	NaN	NaN	0.974466	-2.006747
MGwb4	NaN	NaN	NaN	-0.727707
a60mV	1.431256	1.340309	-1.170299	-0.226169
eW4ge	-0.013960	-0.362543	NaN	NaN
fz64l	0.895717	0.805244	-1.206412	NaN
fzroS	-1.413681	1.607920	1.024180	0.569605

```
nJL3w -0.076467 -1.187678 1.130127 -1.436737
xCmnX      NaN      NaN      NaN -1.219217
```

Note that the row indexes have been unioned and sorted. Here is the same thing with `join='inner'`:

```
In [833]: concat([df.ix[:7, ['a', 'b']], df.ix[2:-2, ['c']],
.....:          df.ix[-7:, ['d']]), axis=1, join='inner')
.....:
Out[833]:
```

	a	b	c	d
a60mV	1.431256	1.340309	-1.170299	-0.226169
0tnT9	0.410835	0.813850	0.132003	-0.827317
nJL3w	-0.076467	-1.187678	1.130127	-1.436737
fzroS	-1.413681	1.607920	1.024180	0.569605

Lastly, suppose we just wanted to reuse the *exact index* from the original DataFrame:

```
In [834]: concat([df.ix[:7, ['a', 'b']], df.ix[2:-2, ['c']],
.....:          df.ix[-7:, ['d']]), axis=1, join_axes=[df.index])
.....:
Out[834]:
```

	a	b	c	d
2dY1o	-1.294524	0.413738	NaN	NaN
eW4ge	-0.013960	-0.362543	NaN	NaN
fz64l	0.895717	0.805244	-1.206412	NaN
a60mV	1.431256	1.340309	-1.170299	-0.226169
0tnT9	0.410835	0.813850	0.132003	-0.827317
nJL3w	-0.076467	-1.187678	1.130127	-1.436737
fzroS	-1.413681	1.607920	1.024180	0.569605
5pSvW	NaN	NaN	0.974466	-2.006747
xCmnX	NaN	NaN	NaN	-1.219217
MGwb4	NaN	NaN	NaN	-0.727707

## 11.1.2 Concatenating using `append`

A useful shortcut to `concat` are the `append` instance methods on `Series` and `DataFrame`. These methods actually predated `concat`. They concatenate along `axis=0`, namely the index:

```
In [835]: s = Series(randn(10), index=np.arange(10))
```

```
In [836]: s1 = s[:5] # note we're slicing with labels here, so 5 is included
```

```
In [837]: s2 = s[6:]
```

```
In [838]: s1.append(s2)
```

```
Out[838]:
0    -0.121306
1    -0.097883
2     0.695775
3     0.341734
4     0.959726
5    -0.619976
6     0.149748
7    -0.732339
8     0.687738
```

In the case of `DataFrame`, the indexes must be disjoint but the columns do not need to be:

```
In [839]: df = DataFrame(randn(6, 4), index=date_range('1/1/2000', periods=6),
.....:                  columns=['A', 'B', 'C', 'D'])
.....:
```

```
In [840]: df1 = df.ix[:3]
```

```
In [841]: df2 = df.ix[3:, :3]
```

```
In [842]: df1
```

```
Out [842]:
```

	A	B	C	D
2000-01-01	0.176444	0.403310	-0.154951	0.301624
2000-01-02	-2.179861	-1.369849	-0.954208	1.462696
2000-01-03	-1.743161	-0.826591	-0.345352	1.314232

```
In [843]: df2
```

```
Out [843]:
```

	A	B	C
2000-01-04	0.690579	0.995761	2.396780
2000-01-05	3.357427	-0.317441	-1.236269
2000-01-06	-0.487602	-0.082240	-2.182937

```
In [844]: df1.append(df2)
```

```
Out [844]:
```

	A	B	C	D
2000-01-01	0.176444	0.403310	-0.154951	0.301624
2000-01-02	-2.179861	-1.369849	-0.954208	1.462696
2000-01-03	-1.743161	-0.826591	-0.345352	1.314232
2000-01-04	0.690579	0.995761	2.396780	NaN
2000-01-05	3.357427	-0.317441	-1.236269	NaN
2000-01-06	-0.487602	-0.082240	-2.182937	NaN

append may take multiple objects to concatenate:

```
In [845]: df1 = df.ix[:2]
```

```
In [846]: df2 = df.ix[2:4]
```

```
In [847]: df3 = df.ix[4:]
```

```
In [848]: df1.append([df2, df3])
```

```
Out [848]:
```

	A	B	C	D
2000-01-01	0.176444	0.403310	-0.154951	0.301624
2000-01-02	-2.179861	-1.369849	-0.954208	1.462696
2000-01-03	-1.743161	-0.826591	-0.345352	1.314232
2000-01-04	0.690579	0.995761	2.396780	0.014871
2000-01-05	3.357427	-0.317441	-1.236269	0.896171
2000-01-06	-0.487602	-0.082240	-2.182937	0.380396

---

**Note:** Unlike *list.append* method, which appends to the original list and returns nothing, append here **does not** modify df1 and returns its copy with df2 appended.

---

### 11.1.3 Ignoring indexes on the concatenation axis

For DataFrames which don't have a meaningful index, you may wish to append them and ignore the fact that they may have overlapping indexes:

```
In [849]: df1 = DataFrame(randn(6, 4), columns=['A', 'B', 'C', 'D'])
```

```
In [850]: df2 = DataFrame(randn(3, 4), columns=['A', 'B', 'C', 'D'])
```

```
In [851]: df1
```

```
Out[851]:
```

	A	B	C	D
0	0.084844	0.432390	1.519970	-0.493662
1	0.600178	0.274230	0.132885	-0.023688
2	2.410179	1.450520	0.206053	-0.251905
3	-2.213588	1.063327	1.266143	0.299368
4	-0.863838	0.408204	-1.048089	-0.025747
5	-0.988387	0.094055	1.262731	1.289997

```
In [852]: df2
```

```
Out[852]:
```

	A	B	C	D
0	0.082423	-0.055758	0.536580	-0.489682
1	0.369374	-0.034571	-2.484478	-0.281461
2	0.030711	0.109121	1.126203	-0.977349

To do this, use the `ignore_index` argument:

```
In [853]: concat([df1, df2], ignore_index=True)
```

```
Out[853]:
```

	A	B	C	D
0	0.084844	0.432390	1.519970	-0.493662
1	0.600178	0.274230	0.132885	-0.023688
2	2.410179	1.450520	0.206053	-0.251905
3	-2.213588	1.063327	1.266143	0.299368
4	-0.863838	0.408204	-1.048089	-0.025747
5	-0.988387	0.094055	1.262731	1.289997
6	0.082423	-0.055758	0.536580	-0.489682
7	0.369374	-0.034571	-2.484478	-0.281461
8	0.030711	0.109121	1.126203	-0.977349

This is also a valid argument to `DataFrame.append`:

```
In [854]: df1.append(df2, ignore_index=True)
```

```
Out[854]:
```

	A	B	C	D
0	0.084844	0.432390	1.519970	-0.493662
1	0.600178	0.274230	0.132885	-0.023688
2	2.410179	1.450520	0.206053	-0.251905
3	-2.213588	1.063327	1.266143	0.299368
4	-0.863838	0.408204	-1.048089	-0.025747
5	-0.988387	0.094055	1.262731	1.289997
6	0.082423	-0.055758	0.536580	-0.489682
7	0.369374	-0.034571	-2.484478	-0.281461
8	0.030711	0.109121	1.126203	-0.977349

### 11.1.4 More concatenating with group keys

Let's consider a variation on the first example presented:

```
In [855]: df = DataFrame(np.random.randn(10, 4))
```

```
In [856]: df
```

```
Out[856]:
```

	0	1	2	3
0	1.474071	-0.064034	-1.282782	0.781836
1	-1.071357	0.441153	2.353925	0.583787
2	0.221471	-0.744471	0.758527	1.729689
3	-0.964980	-0.845696	-1.340896	1.846883
4	-1.328865	1.682706	-1.717693	0.888782
5	0.228440	0.901805	1.171216	0.520260
6	-1.197071	-1.066969	-0.303421	-0.858447
7	0.306996	-0.028665	0.384316	1.574159
8	1.588931	0.476720	0.473424	-0.242861
9	-0.014805	-0.284319	0.650776	-1.461665

```
# break it into pieces
```

```
In [857]: pieces = [df.ix[:, [0, 1]], df.ix[:, [2]], df.ix[:, [3]]]
```

```
In [858]: result = concat(pieces, axis=1, keys=['one', 'two', 'three'])
```

```
In [859]: result
```

```
Out[859]:
```

	one		two		three
	0	1	2	3	
0	1.474071	-0.064034	-1.282782	0.781836	
1	-1.071357	0.441153	2.353925	0.583787	
2	0.221471	-0.744471	0.758527	1.729689	
3	-0.964980	-0.845696	-1.340896	1.846883	
4	-1.328865	1.682706	-1.717693	0.888782	
5	0.228440	0.901805	1.171216	0.520260	
6	-1.197071	-1.066969	-0.303421	-0.858447	
7	0.306996	-0.028665	0.384316	1.574159	
8	1.588931	0.476720	0.473424	-0.242861	
9	-0.014805	-0.284319	0.650776	-1.461665	

You can also pass a dict to `concat` in which case the dict keys will be used for the `keys` argument (unless other keys are specified):

```
In [860]: pieces = {'one': df.ix[:, [0, 1]],
.....:               'two': df.ix[:, [2]],
.....:               'three': df.ix[:, [3]]}
.....:
```

```
In [861]: concat(pieces, axis=1)
```

```
Out[861]:
```

	one		three	two
	0	1	3	2
0	1.474071	-0.064034	0.781836	-1.282782
1	-1.071357	0.441153	0.583787	2.353925
2	0.221471	-0.744471	1.729689	0.758527
3	-0.964980	-0.845696	1.846883	-1.340896
4	-1.328865	1.682706	0.888782	-1.717693
5	0.228440	0.901805	0.520260	1.171216
6	-1.197071	-1.066969	-0.858447	-0.303421

```
7  0.306996 -0.028665  1.574159  0.384316
8  1.588931  0.476720 -0.242861  0.473424
9 -0.014805 -0.284319 -1.461665  0.650776
```

```
In [862]: concat(pieces, keys=['three', 'two'])
```

```
Out[862]:
```

		2	3
three	0	NaN	0.781836
	1	NaN	0.583787
	2	NaN	1.729689
	3	NaN	1.846883
	4	NaN	0.888782
	5	NaN	0.520260
	6	NaN	-0.858447
	7	NaN	1.574159
	8	NaN	-0.242861
	9	NaN	-1.461665
two	0	-1.282782	NaN
	1	2.353925	NaN
	2	0.758527	NaN
	3	-1.340896	NaN
	4	-1.717693	NaN
	5	1.171216	NaN
	6	-0.303421	NaN
	7	0.384316	NaN
	8	0.473424	NaN
	9	0.650776	NaN

The MultiIndex created has levels that are constructed from the passed keys and the columns of the DataFrame pieces:

```
In [863]: result.columns.levels
```

```
Out[863]: [Index([one, two, three], dtype=object), Int64Index([0, 1, 2, 3])]
```

If you wish to specify other levels (as will occasionally be the case), you can do so using the levels argument:

```
In [864]: result = concat(pieces, axis=1, keys=['one', 'two', 'three'],
.....:                    levels=[['three', 'two', 'one', 'zero']],
.....:                    names=['group_key'])
.....:
```

```
In [865]: result
```

```
Out[865]:
```

group_key	one		two	three
	0	1	2	3
0	1.474071	-0.064034	-1.282782	0.781836
1	-1.071357	0.441153	2.353925	0.583787
2	0.221471	-0.744471	0.758527	1.729689
3	-0.964980	-0.845696	-1.340896	1.846883
4	-1.328865	1.682706	-1.717693	0.888782
5	0.228440	0.901805	1.171216	0.520260
6	-1.197071	-1.066969	-0.303421	-0.858447
7	0.306996	-0.028665	0.384316	1.574159
8	1.588931	0.476720	0.473424	-0.242861
9	-0.014805	-0.284319	0.650776	-1.461665

```
In [866]: result.columns.levels
```

```
Out[866]: [Index([three, two, one, zero], dtype=object), Int64Index([0, 1, 2, 3])]
```

Yes, this is fairly esoteric, but is actually necessary for implementing things like GroupBy where the order of a



categorical variable is meaningful.

### 11.1.5 Appending rows to a DataFrame

While not especially efficient (since a new object must be created), you can append a single row to a DataFrame by passing a Series or dict to append, which returns a new DataFrame as above.

```
In [867]: df = DataFrame(np.random.randn(8, 4), columns=['A', 'B', 'C', 'D'])
```

```
In [868]: df
```

```
Out[868]:
```

	A	B	C	D
0	-1.137707	-0.891060	-0.693921	1.613616
1	0.464000	0.227371	-0.496922	0.306389
2	-2.290613	-1.134623	-1.561819	-0.260838
3	0.281957	1.523962	-0.902937	0.068159
4	-0.057873	-0.368204	-1.144073	0.861209
5	0.800193	0.782098	-1.069094	-1.099248
6	0.255269	0.009750	0.661084	0.379319
7	-0.008434	1.952541	-1.056652	0.533946

```
In [869]: s = df.xs(3)
```

```
In [870]: df.append(s, ignore_index=True)
```

```
Out[870]:
```

	A	B	C	D
0	-1.137707	-0.891060	-0.693921	1.613616
1	0.464000	0.227371	-0.496922	0.306389
2	-2.290613	-1.134623	-1.561819	-0.260838
3	0.281957	1.523962	-0.902937	0.068159
4	-0.057873	-0.368204	-1.144073	0.861209
5	0.800193	0.782098	-1.069094	-1.099248
6	0.255269	0.009750	0.661084	0.379319
7	-0.008434	1.952541	-1.056652	0.533946
8	0.281957	1.523962	-0.902937	0.068159

You should use `ignore_index` with this method to instruct DataFrame to discard its index. If you wish to preserve the index, you should construct an appropriately-indexed DataFrame and append or concatenate those objects.

You can also pass a list of dicts or Series:

```
In [871]: df = DataFrame(np.random.randn(5, 4),
.....:                  columns=['foo', 'bar', 'baz', 'qux'])
.....:
```

```
In [872]: dicts = [{'foo': 1, 'bar': 2, 'baz': 3, 'peekaboo': 4},
.....:             {'foo': 5, 'bar': 6, 'baz': 7, 'peekaboo': 8}]
.....:
```

```
In [873]: result = df.append(dicts, ignore_index=True)
```

```
In [874]: result
```

```
Out[874]:
```

	bar	baz	foo	peekaboo	qux
0	0.040403	-0.507516	-1.226970	NaN	-0.230096
1	-1.934370	-1.652499	0.394500	NaN	1.488753
2	0.576897	1.146000	-0.896484	NaN	1.487349
3	2.121453	0.597701	0.604603	NaN	0.563700

```
4 -1.057909  1.375020  0.967661      NaN -0.928797
5  2.000000  3.000000  1.000000         4      NaN
6  6.000000  7.000000  5.000000         8      NaN
```

## 11.2 Database-style DataFrame joining/merging

pandas has full-featured, **high performance** in-memory join operations idiomatically very similar to relational databases like SQL. These methods perform significantly better (in some cases well over an order of magnitude better) than other open source implementations (like `base::merge.data.frame` in R). The reason for this is careful algorithmic design and internal layout of the data in `DataFrame`.

pandas provides a single function, `merge`, as the entry point for all standard database join operations between `DataFrame` objects:

```
merge(left, right, how='left', on=None, left_on=None, right_on=None,
      left_index=False, right_index=False, sort=True,
      suffixes=('.x', '.y'), copy=True)
```

Here's a description of what each argument is for:

- `left`: A `DataFrame` object
- `right`: Another `DataFrame` object
- `on`: Columns (names) to join on. Must be found in both the left and right `DataFrame` objects. If not passed and `left_index` and `right_index` are `False`, the intersection of the columns in the `DataFrames` will be inferred to be the join keys
- `left_on`: Columns from the left `DataFrame` to use as keys. Can either be column names or arrays with length equal to the length of the `DataFrame`
- `right_on`: Columns from the right `DataFrame` to use as keys. Can either be column names or arrays with length equal to the length of the `DataFrame`
- `left_index`: If `True`, use the index (row labels) from the left `DataFrame` as its join key(s). In the case of a `DataFrame` with a `MultiIndex` (hierarchical), the number of levels must match the number of join keys from the right `DataFrame`
- `right_index`: Same usage as `left_index` for the right `DataFrame`
- `how`: One of `'left'`, `'right'`, `'outer'`, `'inner'`. Defaults to `inner`. See below for more detailed description of each method
- `sort`: Sort the result `DataFrame` by the join keys in lexicographical order. Defaults to `True`, setting to `False` will improve performance substantially in many cases
- `suffixes`: A tuple of string suffixes to apply to overlapping columns. Defaults to `('.x', '.y')`.
- `copy`: Always copy data (default `True`) from the passed `DataFrame` objects, even when reindexing is not necessary. Cannot be avoided in many cases but may improve performance / memory usage. The cases where copying can be avoided are somewhat pathological but this option is provided nonetheless.

`merge` is a function in the pandas namespace, and it is also available as a `DataFrame` instance method, with the calling `DataFrame` being implicitly considered the left object in the join.

The related `DataFrame.join` method, uses `merge` internally for the index-on-index and index-on-column(s) joins, but *joins on indexes* by default rather than trying to join on common columns (the default behavior for `merge`). If you are joining on index, you may wish to use `DataFrame.join` to save yourself some typing.

### 11.2.1 Brief primer on merge methods (relational algebra)

Experienced users of relational databases like SQL will be familiar with the terminology used to describe join operations between two SQL-table like structures (DataFrame objects). There are several cases to consider which are very important to understand:

- **one-to-one** joins: for example when joining two DataFrame objects on their indexes (which must contain unique values)
- **many-to-one** joins: for example when joining an index (unique) to one or more columns in a DataFrame
- **many-to-many** joins: joining columns on columns.

---

**Note:** When joining columns on columns (potentially a many-to-many join), any indexes on the passed DataFrame objects **will be discarded**.

---

It is worth spending some time understanding the result of the **many-to-many** join case. In SQL / standard relational algebra, if a key combination appears more than once in both tables, the resulting table will have the **Cartesian product** of the associated data. Here is a very basic example with one unique key combination:

```
In [875]: left = DataFrame({'key': ['foo', 'foo'], 'lval': [1, 2]})
```

```
In [876]: right = DataFrame({'key': ['foo', 'foo'], 'rval': [4, 5]})
```

```
In [877]: left
```

```
Out[877]:
   key  lval
0  foo     1
1  foo     2
```

```
In [878]: right
```

```
Out[878]:
   key  rval
0  foo     4
1  foo     5
```

```
In [879]: merge(left, right, on='key')
```

```
Out[879]:
   key  lval  rval
0  foo     1     4
1  foo     1     5
2  foo     2     4
3  foo     2     5
```

Here is a more complicated example with multiple join keys:

```
In [880]: left = DataFrame({'key1': ['foo', 'foo', 'bar'],
.....:                    'key2': ['one', 'two', 'one'],
.....:                    'lval': [1, 2, 3]})
.....:
```

```
In [881]: right = DataFrame({'key1': ['foo', 'foo', 'bar', 'bar'],
.....:                      'key2': ['one', 'one', 'one', 'two'],
.....:                      'rval': [4, 5, 6, 7]})
.....:
```

```
In [882]: merge(left, right, how='outer')
```

```
Out[882]:
   key1 key2  lval  rval
0  foo  one     1     4
1  foo  one     2     5
2  bar  one     3     6
3  bar  two     0     7
```

```
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
```

```
In [883]: merge(left, right, how='inner')
```

```
Out[883]:
   key1 key2 lval rval
0 bar one 3 6
1 foo one 1 4
2 foo one 1 5
```

The `how` argument to `merge` specifies how to determine which keys are to be included in the resulting table. If a key combination **does not appear** in either the left or right tables, the values in the joined table will be NA. Here is a summary of the `how` options and their SQL equivalent names:

Merge method	SQL Join Name	Description
left	LEFT OUTER JOIN	Use keys from left frame only
right	RIGHT OUTER JOIN	Use keys from right frame only
outer	FULL OUTER JOIN	Use union of keys from both frames
inner	INNER JOIN	Use intersection of keys from both frames

Note that if using the index from either the left or right DataFrame (or both) using the `left_index/ right_index` options, the join operation is no longer a many-to-many join by construction, as the index values are necessarily unique. There will be some examples of this below.

## 11.2.2 Joining on index

`DataFrame.join` is a convenient method for combining the columns of two potentially differently-indexed DataFrames into a single result DataFrame. Here is a very basic example:

```
In [884]: df = DataFrame(np.random.randn(8, 4), columns=['A', 'B', 'C', 'D'])
```

```
In [885]: df1 = df.ix[1:, ['A', 'B']]
```

```
In [886]: df2 = df.ix[:5, ['C', 'D']]
```

```
In [887]: df1
```

```
Out[887]:
      A      B
1 -2.461467 -1.553902
2  1.771740 -0.670027
3 -3.201750  0.792716
4 -0.747169 -0.309038
5  0.936527  1.255746
6  0.062297 -0.110388
7  0.077849  0.629498
```

```
In [888]: df2
```

```
Out[888]:
      C      D
0  0.377953  0.493672
1  2.015523 -1.833722
2  0.049307 -0.521493
3  0.146111  1.903247
4  0.393876  1.861468
```

```
5 -2.655452  1.219492
```

```
In [889]: df1.join(df2)
```

```
Out[889]:
```

	A	B	C	D
1	-2.461467	-1.553902	2.015523	-1.833722
2	1.771740	-0.670027	0.049307	-0.521493
3	-3.201750	0.792716	0.146111	1.903247
4	-0.747169	-0.309038	0.393876	1.861468
5	0.936527	1.255746	-2.655452	1.219492
6	0.062297	-0.110388	NaN	NaN
7	0.077849	0.629498	NaN	NaN

```
In [890]: df1.join(df2, how='outer')
```

```
Out[890]:
```

	A	B	C	D
0	NaN	NaN	0.377953	0.493672
1	-2.461467	-1.553902	2.015523	-1.833722
2	1.771740	-0.670027	0.049307	-0.521493
3	-3.201750	0.792716	0.146111	1.903247
4	-0.747169	-0.309038	0.393876	1.861468
5	0.936527	1.255746	-2.655452	1.219492
6	0.062297	-0.110388	NaN	NaN
7	0.077849	0.629498	NaN	NaN

```
In [891]: df1.join(df2, how='inner')
```

```
Out[891]:
```

	A	B	C	D
1	-2.461467	-1.553902	2.015523	-1.833722
2	1.771740	-0.670027	0.049307	-0.521493
3	-3.201750	0.792716	0.146111	1.903247
4	-0.747169	-0.309038	0.393876	1.861468
5	0.936527	1.255746	-2.655452	1.219492

The data alignment here is on the indexes (row labels). This same behavior can be achieved using `merge` plus additional arguments instructing it to use the indexes:

```
In [892]: merge(df1, df2, left_index=True, right_index=True, how='outer')
```

```
Out[892]:
```

	A	B	C	D
0	NaN	NaN	0.377953	0.493672
1	-2.461467	-1.553902	2.015523	-1.833722
2	1.771740	-0.670027	0.049307	-0.521493
3	-3.201750	0.792716	0.146111	1.903247
4	-0.747169	-0.309038	0.393876	1.861468
5	0.936527	1.255746	-2.655452	1.219492
6	0.062297	-0.110388	NaN	NaN
7	0.077849	0.629498	NaN	NaN

### 11.2.3 Joining key columns on an index

`join` takes an optional `on` argument which may be a column or multiple column names, which specifies that the passed DataFrame is to be aligned on that column in the DataFrame. These two function calls are completely equivalent:

```
left.join(right, on=key_or_keys)
merge(left, right, left_on=key_or_keys, right_index=True,
```

```
how='left', sort=False)
```

Obviously you can choose whichever form you find more convenient. For many-to-one joins (where one of the DataFrame's is already indexed by the join key), using `join` may be more convenient. Here is a simple example:

```
In [893]: df['key'] = ['foo', 'bar'] * 4
```

```
In [894]: to_join = DataFrame(randn(2, 2), index=['bar', 'foo'],
.....:                        columns=['j1', 'j2'])
.....:
```

```
In [895]: df
```

```
Out[895]:
```

	A	B	C	D	key
0	-0.308853	-0.681087	0.377953	0.493672	foo
1	-2.461467	-1.553902	2.015523	-1.833722	bar
2	1.771740	-0.670027	0.049307	-0.521493	foo
3	-3.201750	0.792716	0.146111	1.903247	bar
4	-0.747169	-0.309038	0.393876	1.861468	foo
5	0.936527	1.255746	-2.655452	1.219492	bar
6	0.062297	-0.110388	-1.184357	-0.558081	foo
7	0.077849	0.629498	-1.035260	-0.438229	bar

```
In [896]: to_join
```

```
Out[896]:
```

	j1	j2
bar	0.503703	0.413086
foo	-1.139050	0.660342

```
In [897]: df.join(to_join, on='key')
```

```
Out[897]:
```

	A	B	C	D	key	j1	j2
0	-0.308853	-0.681087	0.377953	0.493672	foo	-1.139050	0.660342
1	-2.461467	-1.553902	2.015523	-1.833722	bar	0.503703	0.413086
2	1.771740	-0.670027	0.049307	-0.521493	foo	-1.139050	0.660342
3	-3.201750	0.792716	0.146111	1.903247	bar	0.503703	0.413086
4	-0.747169	-0.309038	0.393876	1.861468	foo	-1.139050	0.660342
5	0.936527	1.255746	-2.655452	1.219492	bar	0.503703	0.413086
6	0.062297	-0.110388	-1.184357	-0.558081	foo	-1.139050	0.660342
7	0.077849	0.629498	-1.035260	-0.438229	bar	0.503703	0.413086

```
In [898]: merge(df, to_join, left_on='key', right_index=True,
```

```
.....:          how='left', sort=False)
```

```
.....:
```

```
Out[898]:
```

	A	B	C	D	key	j1	j2
0	-0.308853	-0.681087	0.377953	0.493672	foo	-1.139050	0.660342
1	-2.461467	-1.553902	2.015523	-1.833722	bar	0.503703	0.413086
2	1.771740	-0.670027	0.049307	-0.521493	foo	-1.139050	0.660342
3	-3.201750	0.792716	0.146111	1.903247	bar	0.503703	0.413086
4	-0.747169	-0.309038	0.393876	1.861468	foo	-1.139050	0.660342
5	0.936527	1.255746	-2.655452	1.219492	bar	0.503703	0.413086
6	0.062297	-0.110388	-1.184357	-0.558081	foo	-1.139050	0.660342
7	0.077849	0.629498	-1.035260	-0.438229	bar	0.503703	0.413086

To join on multiple keys, the passed DataFrame must have a MultiIndex:

```
In [899]: index = MultiIndex(levels=[['foo', 'bar', 'baz', 'qux'],
.....:                               ['one', 'two', 'three']],
```

```

.....:             labels=[0, 0, 0, 1, 1, 2, 2, 3, 3, 3],
.....:                     [0, 1, 2, 0, 1, 1, 2, 0, 1, 2]],
.....:             names=['first', 'second'])
.....:

In [900]: to_join = DataFrame(np.random.randn(10, 3), index=index,
.....:                        columns=['j_one', 'j_two', 'j_three'])
.....:

# a little relevant example with NAs
In [901]: key1 = ['bar', 'bar', 'bar', 'foo', 'foo', 'baz', 'baz', 'qux',
.....:            'qux', 'snap']
.....:

In [902]: key2 = ['two', 'one', 'three', 'one', 'two', 'one', 'two', 'two',
.....:            'three', 'one']
.....:

In [903]: data = np.random.randn(len(key1))

In [904]: data = DataFrame({'key1' : key1, 'key2' : key2,
.....:                     'data' : data})
.....:

In [905]: data
Out[905]:
   data  key1  key2
0 -1.004168  bar   two
1 -1.377627  bar   one
2  0.499281  bar three
3 -1.405256  foo   one
4  0.162565  foo   two
5 -0.067785  baz   one
6 -1.260006  baz   two
7 -1.132896  qux   two
8 -2.006481  qux three
9  0.301016  snap  one

In [906]: to_join
Out[906]:
           j_one  j_two  j_three
first second
foo   one    0.464794 -0.309337 -0.649593
      two    0.683758 -0.643834  0.421287
      three  1.032814 -1.290493  0.787872
bar   one    1.515707 -0.276487 -0.223762
      two    1.397431  1.503874 -0.478905
baz   two   -0.135950 -0.730327 -0.033277
      three  0.281151 -1.298915 -2.819487
qux   one   -0.851985 -1.106952 -0.937731
      two   -1.537770  0.555759 -2.277282
      three -0.390201  1.207122  0.178690

```

Now this can be joined by passing the two key column names:

```

In [907]: data.join(to_join, on=['key1', 'key2'])
Out[907]:
   data  key1  key2  j_one  j_two  j_three
0 -1.004168  bar   two  1.397431  1.503874 -0.478905

```

```
1 -1.377627 bar one 1.515707 -0.276487 -0.223762
2 0.499281 bar three NaN NaN NaN
3 -1.405256 foo one 0.464794 -0.309337 -0.649593
4 0.162565 foo two 0.683758 -0.643834 0.421287
5 -0.067785 baz one NaN NaN NaN
6 -1.260006 baz two -0.135950 -0.730327 -0.033277
7 -1.132896 qux two -1.537770 0.555759 -2.277282
8 -2.006481 qux three -0.390201 1.207122 0.178690
9 0.301016 snap one NaN NaN NaN
```

The default for `DataFrame.join` is to perform a left join (essentially a “VLOOKUP” operation, for Excel users), which uses only the keys found in the calling `DataFrame`. Other join types, for example inner join, can be just as easily performed:

```
In [908]: data.join(to_join, on=['key1', 'key2'], how='inner')
Out[908]:
```

```
      data key1 key2 j_one j_two j_three
0 -1.004168 bar two 1.397431 1.503874 -0.478905
1 -1.377627 bar one 1.515707 -0.276487 -0.223762
3 -1.405256 foo one 0.464794 -0.309337 -0.649593
4 0.162565 foo two 0.683758 -0.643834 0.421287
6 -1.260006 baz two -0.135950 -0.730327 -0.033277
7 -1.132896 qux two -1.537770 0.555759 -2.277282
8 -2.006481 qux three -0.390201 1.207122 0.178690
```

As you can see, this drops any rows where there was no match.

## 11.2.4 Overlapping value columns

The merge `suffixes` argument takes a tuple of list of strings to append to overlapping column names in the input `DataFrames` to disambiguate the result columns:

```
In [909]: left = DataFrame({'key': ['foo', 'foo'], 'value': [1, 2]})
```

```
In [910]: right = DataFrame({'key': ['foo', 'foo'], 'value': [4, 5]})
```

```
In [911]: merge(left, right, on='key', suffixes=['_left', '_right'])
```

```
Out[911]:
   key  value_left  value_right
0  foo           1            4
1  foo           1            5
2  foo           2            4
3  foo           2            5
```

`DataFrame.join` has `lsuffix` and `rsuffix` arguments which behave similarly.

## 11.2.5 Merging Ordered Data

New in v0.8.0 is the `ordered_merge` function for combining time series and other ordered data. In particular it has an optional `fill_method` keyword to fill/interpolate missing data:

```
In [912]: A
```

```
Out[912]:
   group key  lvalue
0      a   a       1
1      a   c       2
```



```

2      a  e      3
3      b  a      1
4      b  c      2
5      b  e      3

```

```
In [913]: B
```

```
Out[913]:
   key  rvalue
0    b        1
1    c        2
2    d        3

```

```
In [914]: ordered_merge(A, B, fill_method='ffill', left_by='group')
```

```
Out[914]:
   group key  lvalue  rvalue
0      a  a        1     NaN
1      a  b        1        1
2      a  c        2        2
3      a  d        2        3
4      a  e        3        3
5      b  a        1     NaN
6      b  b        1        1
7      b  c        2        2
8      b  d        2        3
9      b  e        3        3

```

## 11.2.6 Joining multiple DataFrame or Panel objects

A list or tuple of DataFrames can also be passed to `DataFrame.join` to join them together on their indexes. The same is true for `Panel.join`.

```
In [915]: df1 = df.ix[:, ['A', 'B']]
```

```
In [916]: df2 = df.ix[:, ['C', 'D']]
```

```
In [917]: df3 = df.ix[:, ['key']]
```

```
In [918]: df1
```

```
Out[918]:
      A      B
0 -0.308853 -0.681087
1 -2.461467 -1.553902
2  1.771740 -0.670027
3 -3.201750  0.792716
4 -0.747169 -0.309038
5  0.936527  1.255746
6  0.062297 -0.110388
7  0.077849  0.629498

```

```
In [919]: df1.join([df2, df3])
```

```
Out[919]:
      A      B      C      D  key
0 -0.308853 -0.681087  0.377953  0.493672  foo
1 -2.461467 -1.553902  2.015523 -1.833722  bar
2  1.771740 -0.670027  0.049307 -0.521493  foo
3 -3.201750  0.792716  0.146111  1.903247  bar
4 -0.747169 -0.309038  0.393876  1.861468  foo

```

```
5  0.936527  1.255746 -2.655452  1.219492  bar
6  0.062297 -0.110388 -1.184357 -0.558081  foo
7  0.077849  0.629498 -1.035260 -0.438229  bar
```

### 11.2.7 Merging together values within Series or DataFrame columns

Another fairly common situation is to have two like-indexed (or similarly indexed) Series or DataFrame objects and wanting to “patch” values in one object from values for matching indices in the other. Here is an example:

```
In [920]: df1 = DataFrame([[nan, 3., 5.], [-4.6, np.nan, nan],
.....:                  [nan, 7., nan]])
.....:

In [921]: df2 = DataFrame([[-42.6, np.nan, -8.2], [-5., 1.6, 4]],
.....:                  index=[1, 2])
.....:
```

For this, use the `combine_first` method:

```
In [922]: df1.combine_first(df2)
Out[922]:
   0    1    2
0  NaN    3  5.0
1 -4.6  NaN -8.2
2 -5.0    7  4.0
```

Note that this method only takes values from the right DataFrame if they are missing in the left DataFrame. A related method, `update`, alters non-NA values inplace:

```
In [923]: df1.update(df2)

In [924]: df1
Out[924]:
   0    1    2
0  NaN    3  5.0
1 -42.6  NaN -8.2
2 -5.0  1.6  4.0
```

# RESHAPING AND PIVOT TABLES

## 12.1 Reshaping by pivoting DataFrame objects

Data is often stored in CSV files or databases in so-called “stacked” or “record” format:

```
In [993]: df
Out[993]:
```

	date	variable	value
0	2000-01-03 00:00:00	A	0.469112
1	2000-01-04 00:00:00	A	-0.282863
2	2000-01-05 00:00:00	A	-1.509059
3	2000-01-03 00:00:00	B	-1.135632
4	2000-01-04 00:00:00	B	1.212112
5	2000-01-05 00:00:00	B	-0.173215
6	2000-01-03 00:00:00	C	0.119209
7	2000-01-04 00:00:00	C	-1.044236
8	2000-01-05 00:00:00	C	-0.861849
9	2000-01-03 00:00:00	D	-2.104569
10	2000-01-04 00:00:00	D	-0.494929
11	2000-01-05 00:00:00	D	1.071804

For the curious here is how the above DataFrame was created:

```
import pandas.util.testing as tm; tm.N = 3
def unpivot(frame):
    N, K = frame.shape
    data = {'value' : frame.values.ravel('F'),
           'variable' : np.asarray(frame.columns).repeat(N),
           'date' : np.tile(np.asarray(frame.index), K)}
    return DataFrame(data, columns=['date', 'variable', 'value'])
df = unpivot(tm.makeTimeDataFrame())
```

To select out everything for variable A we could do:

```
In [994]: df[df['variable'] == 'A']
Out[994]:
```

	date	variable	value
0	2000-01-03 00:00:00	A	0.469112
1	2000-01-04 00:00:00	A	-0.282863
2	2000-01-05 00:00:00	A	-1.509059

But suppose we wish to do time series operations with the variables. A better representation would be where the columns are the unique variables and an index of dates identifies individual observations. To reshape the data into this form, use the `pivot` function:

```
In [995]: df.pivot(index='date', columns='variable', values='value')
Out[995]:
variable          A          B          C          D
date
2000-01-03  0.469112 -1.135632  0.119209 -2.104569
2000-01-04 -0.282863  1.212112 -1.044236 -0.494929
2000-01-05 -1.509059 -0.173215 -0.861849  1.071804
```

If the `values` argument is omitted, and the input `DataFrame` has more than one column of values which are not used as column or index inputs to `pivot`, then the resulting “pivoted” `DataFrame` will have *hierarchical columns* whose toplevel level indicates the respective value column:

```
In [996]: df['value2'] = df['value'] * 2

In [997]: pivoted = df.pivot('date', 'variable')

In [998]: pivoted
Out[998]:
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 3 entries, 2000-01-03 00:00:00 to 2000-01-05 00:00:00
Data columns:
('value', 'A')      3  non-null values
('value', 'B')      3  non-null values
('value', 'C')      3  non-null values
('value', 'D')      3  non-null values
('value2', 'A')     3  non-null values
('value2', 'B')     3  non-null values
('value2', 'C')     3  non-null values
('value2', 'D')     3  non-null values
dtypes: float64(8)
```

You of course can then select subsets from the pivoted `DataFrame`:

```
In [999]: pivoted['value2']
Out[999]:
variable          A          B          C          D
date
2000-01-03  0.938225 -2.271265  0.238417 -4.209138
2000-01-04 -0.565727  2.424224 -2.088472 -0.989859
2000-01-05 -3.018117 -0.346429 -1.723698  2.143608
```

Note that this returns a view on the underlying data in the case where the data are homogeneously-typed.

## 12.2 Reshaping by stacking and unstacking

Closely related to the `pivot` function are the related `stack` and `unstack` functions currently available on `Series` and `DataFrame`. These functions are designed to work together with `MultiIndex` objects (see the section on *hierarchical indexing*). Here are essentially what these functions do:

- `stack`: “pivot” a level of the (possibly hierarchical) column labels, returning a `DataFrame` with an index with a new inner-most level of row labels.
- `unstack`: inverse operation from `stack`: “pivot” a level of the (possibly hierarchical) row index to the column axis, producing a reshaped `DataFrame` with a new inner-most level of column labels.

The clearest way to explain is by example. Let’s take a prior example data set from the hierarchical indexing section:

```
In [1000]: tuples = zip(*[['bar', 'bar', 'baz', 'baz',
.....:                    'foo', 'foo', 'qux', 'qux'],
.....:                    ['one', 'two', 'one', 'two',
.....:                    'one', 'two', 'one', 'two']])
.....:

In [1001]: index = MultiIndex.from_tuples(tuples, names=['first', 'second'])

In [1002]: df = DataFrame(randn(8, 2), index=index, columns=['A', 'B'])

In [1003]: df2 = df[:4]

In [1004]: df2
Out[1004]:
```

		A	B
first	second		
bar	one	0.721555	-0.706771
	two	-1.039575	0.271860
baz	one	-0.424972	0.567020
	two	0.276232	-1.087401

The stack function “compresses” a level in the DataFrame’s columns to produce either:

- A Series, in the case of a simple column Index
- A DataFrame, in the case of a MultiIndex in the columns

If the columns have a MultiIndex, you can choose which level to stack. The stacked level becomes the new lowest level in a MultiIndex on the columns:

```
In [1005]: stacked = df2.stack()

In [1006]: stacked
Out[1006]:
```

first	second		
bar	one	A	0.721555
		B	-0.706771
	two	A	-1.039575
		B	0.271860
baz	one	A	-0.424972
		B	0.567020
	two	A	0.276232
		B	-1.087401

With a “stacked” DataFrame or Series (having a MultiIndex as the index), the inverse operation of stack is unstack, which by default unstacks the last level:

```
In [1007]: stacked.unstack()
Out[1007]:
```

		A	B
first	second		
bar	one	0.721555	-0.706771
	two	-1.039575	0.271860
baz	one	-0.424972	0.567020
	two	0.276232	-1.087401

```
In [1008]: stacked.unstack(1)
Out[1008]:
```

second	one	two
first		

```
bar   A   0.721555 -1.039575
      B  -0.706771  0.271860
baz   A  -0.424972  0.276232
      B   0.567020 -1.087401
```

```
In [1009]: stacked.unstack(0)
```

```
Out[1009]:
first      bar      baz
second
one   A   0.721555 -0.424972
      B  -0.706771  0.567020
two   A  -1.039575  0.276232
      B   0.271860 -1.087401
```

If the indexes have names, you can use the level names instead of specifying the level numbers:

```
In [1010]: stacked.unstack('second')
```

```
Out[1010]:
second      one      two
first
bar   A   0.721555 -1.039575
      B  -0.706771  0.271860
baz   A  -0.424972  0.276232
      B   0.567020 -1.087401
```

You may also stack or unstack more than one level at a time by passing a list of levels, in which case the end result is as if each level in the list were processed individually.

These functions are intelligent about handling missing data and do not expect each subgroup within the hierarchical index to have the same set of labels. They also can handle the index being unsorted (but you can make it sorted by calling `sortlevel`, of course). Here is a more complex example:

```
In [1011]: columns = MultiIndex.from_tuples([('A', 'cat'), ('B', 'dog'),
.....:                                     ('B', 'cat'), ('A', 'dog')],
.....:                                     names=['exp', 'animal'])
.....:
```

```
In [1012]: df = DataFrame(randn(8, 4), index=index, columns=columns)
```

```
In [1013]: df2 = df.ix[[0, 1, 2, 4, 5, 7]]
```

```
In [1014]: df2
```

```
Out[1014]:
exp      A      B      A
animal   cat   dog   cat   dog
first second
bar   one  -0.370647 -1.157892 -1.344312  0.844885
      two   1.075770 -0.109050  1.643563 -1.469388
baz   one   0.357021 -0.674600 -1.776904 -0.968914
foo   one  -0.013960 -0.362543 -0.006154 -0.923061
      two   0.895717  0.805244 -1.206412  2.565646
qux   two   0.410835  0.813850  0.132003 -0.827317
```

As mentioned above, `stack` can be called with a `level` argument to select which level in the columns to stack:

```
In [1015]: df2.stack('exp')
```

```
Out[1015]:
animal      cat      dog
first second exp
```

```

bar   one    A   -0.370647  0.844885
      two    A    1.075770 -1.469388
      two    B    1.643563 -0.109050
baz   one    A    0.357021 -0.968914
      two    B   -1.776904 -0.674600
foo   one    A   -0.013960 -0.923061
      two    B   -0.006154 -0.362543
qux   two    A    0.895717  2.565646
      two    B   -1.206412  0.805244
      two    A    0.410835 -0.827317
      two    B    0.132003  0.813850

```

```
In [1016]: df2.stack('animal')
```

```
Out[1016]:
```

```

exp          A          B
first second animal
bar   one   cat   -0.370647 -1.344312
      two   dog    0.844885 -1.157892
      two   cat    1.075770  1.643563
      two   dog   -1.469388 -0.109050
baz   one   cat    0.357021 -1.776904
      two   dog   -0.968914 -0.674600
foo   one   cat   -0.013960 -0.006154
      two   dog   -0.923061 -0.362543
      two   cat    0.895717 -1.206412
      two   dog    2.565646  0.805244
qux   two   cat    0.410835  0.132003
      two   dog   -0.827317  0.813850

```

Unstacking when the columns are a MultiIndex is also careful about doing the right thing:

```
In [1017]: df[:3].unstack(0)
```

```
Out[1017]:
```

```

<class 'pandas.core.frame.DataFrame'>
Index: 2 entries, one to two
Data columns:
('A', 'cat', 'bar')    2 non-null values
('A', 'cat', 'baz')    1 non-null values
('B', 'dog', 'bar')    2 non-null values
('B', 'dog', 'baz')    1 non-null values
('B', 'cat', 'bar')    2 non-null values
('B', 'cat', 'baz')    1 non-null values
('A', 'dog', 'bar')    2 non-null values
('A', 'dog', 'baz')    1 non-null values
dtypes: float64(8)

```

```
In [1018]: df2.unstack(1)
```

```
Out[1018]:
```

```

<class 'pandas.core.frame.DataFrame'>
Index: 4 entries, bar to qux
Data columns:
('A', 'cat', 'one')    3 non-null values
('A', 'cat', 'two')    3 non-null values
('B', 'dog', 'one')    3 non-null values
('B', 'dog', 'two')    3 non-null values
('B', 'cat', 'one')    3 non-null values
('B', 'cat', 'two')    3 non-null values

```

```
('A', 'dog', 'one')    3 non-null values
('A', 'dog', 'two')    3 non-null values
dtypes: float64(8)
```

## 12.3 Reshaping by Melt

The `melt` function found in `pandas.core.reshape` is useful to massage a `DataFrame` into a format where one or more columns are identifier variables, while all other columns, considered measured variables, are “pivoted” to the row axis, leaving just two non-identifier columns, “variable” and “value”.

For instance,

```
In [1019]: cheese = DataFrame({'first' : ['John', 'Mary'],
.....:                        'last'  : ['Doe', 'Bo'],
.....:                        'height' : [5.5, 6.0],
.....:                        'weight' : [130, 150]})
.....:
```

```
In [1020]: cheese
```

```
Out[1020]:
```

	first	height	last	weight
0	John	5.5	Doe	130
1	Mary	6.0	Bo	150

```
In [1021]: melt(cheese, id_vars=['first', 'last'])
```

```
Out[1021]:
```

	first	last	variable	value
0	John	Doe	height	5.5
1	Mary	Bo	height	6.0
2	John	Doe	weight	130.0
3	Mary	Bo	weight	150.0

## 12.4 Combining with stats and GroupBy

It should be no shock that combining `pivot / stack / unstack` with `GroupBy` and the basic `Series` and `DataFrame` statistical functions can produce some very expressive and fast data manipulations.

```
In [1022]: df
```

```
Out[1022]:
```

exp		A	B		A
animal		cat	dog	cat	dog
first	second				
bar	one	-0.370647	-1.157892	-1.344312	0.844885
	two	1.075770	-0.109050	1.643563	-1.469388
baz	one	0.357021	-0.674600	-1.776904	-0.968914
	two	-1.294524	0.413738	0.276662	-0.472035
foo	one	-0.013960	-0.362543	-0.006154	-0.923061
	two	0.895717	0.805244	-1.206412	2.565646
qux	one	1.431256	1.340309	-1.170299	-0.226169
	two	0.410835	0.813850	0.132003	-0.827317

```
In [1023]: df.stack().mean(1).unstack()
```

```
Out[1023]:
```

animal		cat	dog
--------	--	-----	-----



```

first second
bar  one  -0.857479 -0.156504
    two   1.359666 -0.789219
baz  one  -0.709942 -0.821757
    two  -0.508931 -0.029148
foo  one  -0.010057 -0.642802
    two  -0.155347  1.685445
qux  one   0.130479  0.557070
    two   0.271419 -0.006733

# same result, another way
In [1024]: df.groupby(level=1, axis=1).mean()
Out[1024]:
animal      cat      dog
first second
bar  one  -0.857479 -0.156504
    two   1.359666 -0.789219
baz  one  -0.709942 -0.821757
    two  -0.508931 -0.029148
foo  one  -0.010057 -0.642802
    two  -0.155347  1.685445
qux  one   0.130479  0.557070
    two   0.271419 -0.006733

In [1025]: df.stack().groupby(level=1).mean()
Out[1025]:
exp      A      B
second
one      0.016301 -0.644049
two      0.110588  0.346200

In [1026]: df.mean().unstack(0)
Out[1026]:
exp      A      B
animal
cat      0.311433 -0.431481
dog     -0.184544  0.133632

```

## 12.5 Pivot tables and cross-tabulations

The function `pandas.pivot_table` can be used to create spreadsheet-style pivot tables. It takes a number of arguments

- `data`: A `DataFrame` object
- `values`: a column or a list of columns to aggregate
- `rows`: list of columns to group by on the table rows
- `cols`: list of columns to group by on the table columns
- `aggfunc`: function to use for aggregation, defaulting to `numpy.mean`

Consider a data set like this:

```

In [1027]: df = DataFrame({'A' : ['one', 'one', 'two', 'three'] * 6,
.....:                    'B' : ['A', 'B', 'C'] * 8,
.....:                    'C' : ['foo', 'foo', 'foo', 'bar', 'bar', 'bar'] * 4,

```

```
.....:         'D' : np.random.randn(24),
.....:         'E' : np.random.randn(24) })
.....:
```

```
In [1028]: df
```

```
Out[1028]:
```

	A	B	C	D	E
0	one	A	foo	-0.076467	0.959726
1	one	B	foo	-1.187678	-1.110336
2	two	C	foo	1.130127	-0.619976
3	three	A	bar	-1.436737	0.149748
4	one	B	bar	-1.413681	-0.732339
5	one	C	bar	1.607920	0.687738
6	two	A	foo	1.024180	0.176444
7	three	B	foo	0.569605	0.403310
8	one	C	foo	0.875906	-0.154951
9	one	A	bar	-2.211372	0.301624
10	two	B	bar	0.974466	-2.179861
11	three	C	bar	-2.006747	-1.369849
12	one	A	foo	-0.410001	-0.954208
13	one	B	foo	-0.078638	1.462696
14	two	C	foo	0.545952	-1.743161
15	three	A	bar	-1.219217	-0.826591
16	one	B	bar	-1.226825	-0.345352
17	one	C	bar	0.769804	1.314232
18	two	A	foo	-1.281247	0.690579
19	three	B	foo	-0.727707	0.995761
20	one	C	foo	-0.121306	2.396780
21	one	A	bar	-0.097883	0.014871
22	two	B	bar	0.695775	3.357427
23	three	C	bar	0.341734	-0.317441

We can produce pivot tables from this data very easily:

```
In [1029]: pivot_table(df, values='D', rows=['A', 'B'], cols=['C'])
```

```
Out[1029]:
```

C	bar	foo
A	B	
one	A	-1.154627 -0.243234
	B	-1.320253 -0.633158
	C	1.188862 0.377300
three	A	-1.327977 NaN
	B	NaN -0.079051
	C	-0.832506 NaN
two	A	NaN -0.128534
	B	0.835120 NaN
	C	NaN 0.838040

```
In [1030]: pivot_table(df, values='D', rows=['B'], cols=['A', 'C'], aggfunc=np.sum)
```

```
Out[1030]:
```

A	one	three	two	
C	bar	foo	bar	foo
B				
A	-2.309255 -0.486468 -2.655954	NaN	NaN	-0.257067
B	-2.640506 -1.266315	NaN -0.158102	1.670241	NaN
C	2.377724 0.754600 -1.665013	NaN	NaN	1.676079

```
In [1031]: pivot_table(df, values=['D', 'E'], rows=['B'], cols=['A', 'C'], aggfunc=np.sum)
```

```
Out[1031]:
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 3 entries, A to C
Data columns:
('D', 'one', 'bar')      3 non-null values
('D', 'one', 'foo')      3 non-null values
('D', 'three', 'bar')    2 non-null values
('D', 'three', 'foo')    1 non-null values
('D', 'two', 'bar')      1 non-null values
('D', 'two', 'foo')      2 non-null values
('E', 'one', 'bar')      3 non-null values
('E', 'one', 'foo')      3 non-null values
('E', 'three', 'bar')    2 non-null values
('E', 'three', 'foo')    1 non-null values
('E', 'two', 'bar')      1 non-null values
('E', 'two', 'foo')      2 non-null values
dtypes: float64(12)
```

The result object is a DataFrame having potentially hierarchical indexes on the rows and columns. If the values column name is not given, the pivot table will include all of the data that can be aggregated in an additional level of hierarchy in the columns:

```
In [1032]: pivot_table(df, rows=['A', 'B'], cols=['C'])
Out[1032]:
```

		D		E	
C		bar	foo	bar	foo
A	one	A	-1.154627 -0.243234	0.158248	0.002759
		B	-1.320253 -0.633158	-0.538846	0.176180
		C	1.188862 0.377300	1.000985	1.120915
three	A	-1.327977	NaN	-0.338421	NaN
	B	NaN	-0.079051	NaN	0.699535
	C	-0.832506	NaN	-0.843645	NaN
two	A	NaN	-0.128534	NaN	0.433512
	B	0.835120	NaN	0.588783	NaN
	C	NaN	0.838040	NaN	-1.181568

You can render a nice output of the table omitting the missing values by calling `to_string` if you wish:

```
In [1033]: table = pivot_table(df, rows=['A', 'B'], cols=['C'])
```

```
In [1034]: print table.to_string(na_rep='')
```

		D		E	
C		bar	foo	bar	foo
A	one	A	-1.154627 -0.243234	0.158248	0.002759
		B	-1.320253 -0.633158	-0.538846	0.176180
		C	1.188862 0.377300	1.000985	1.120915
three	A	-1.327977		-0.338421	
	B		-0.079051		0.699535
	C	-0.832506		-0.843645	
two	A		-0.128534		0.433512
	B	0.835120		0.588783	
	C		0.838040		-1.181568

Note that `pivot_table` is also available as an instance method on `DataFrame`.

### 12.5.1 Cross tabulations

Use the `crosstab` function to compute a cross-tabulation of two (or more) factors. By default `crosstab` computes a frequency table of the factors unless an array of values and an aggregation function are passed.

It takes a number of arguments

- `rows`: array-like, values to group by in the rows
- `cols`: array-like, values to group by in the columns
- `values`: array-like, optional, array of values to aggregate according to the factors
- `aggfunc`: function, optional, If no values array is passed, computes a frequency table
- `rownames`: sequence, default None, must match number of row arrays passed
- `colnames`: sequence, default None, if passed, must match number of column arrays passed
- `margins`: boolean, default False, Add row/column margins (subtotals)

Any Series passed will have their name attributes used unless row or column names for the cross-tabulation are specified

For example:

```
In [1035]: foo, bar, dull, shiny, one, two = 'foo', 'bar', 'dull', 'shiny', 'one', 'two'
```

```
In [1036]: a = np.array([foo, foo, bar, bar, foo, foo], dtype=object)
```

```
In [1037]: b = np.array([one, one, two, one, two, one], dtype=object)
```

```
In [1038]: c = np.array([dull, dull, shiny, dull, dull, shiny], dtype=object)
```

```
In [1039]: crosstab(a, [b, c], rownames=['a'], colnames=['b', 'c'])
```

```
Out[1039]:
b      one      two
c  dull shiny dull shiny
a
bar    1      0      0      1
foo    2      1      1      0
```

### 12.5.2 Adding margins (partial aggregates)

If you pass `margins=True` to `pivot_table`, special All columns and rows will be added with partial group aggregates across the categories on the rows and columns:

```
In [1040]: df.pivot_table(rows=['A', 'B'], cols='C', margins=True, aggfunc=np.std)
```

```
Out[1040]:
C      D      E
      bar  foo  All  bar  foo  All
A  B
one  A  1.494463  0.235844  1.019752  0.202765  1.353355  0.795165
     B  0.132127  0.784210  0.606779  0.273641  1.819408  1.139647
     C  0.592638  0.705136  0.708771  0.442998  1.804346  1.074910
three A  0.153810      NaN  0.153810  0.690376      NaN  0.690376
     B      NaN  0.917338  0.917338      NaN  0.418926  0.418926
     C  1.660627      NaN  1.660627  0.744165      NaN  0.744165
two  A      NaN  1.630183  1.630183      NaN  0.363548  0.363548
     B  0.197065      NaN  0.197065  3.915454      NaN  3.915454
```

```
      C      NaN  0.413074  0.413074      NaN  0.794212  0.794212
All    1.294620  0.824989  1.064129  1.403041  1.188419  1.248988
```

## 12.6 Tiling

The `cut` function computes groupings for the values of the input array and is often used to transform continuous variables to discrete or categorical variables:

```
In [1041]: ages = np.array([10, 15, 13, 12, 23, 25, 28, 59, 60])
```

```
In [1042]: cut(ages, bins=3)
```

```
Out[1042]:
```

```
Categorical:
```

```
array([(9.95, 26.667], (9.95, 26.667], (9.95, 26.667], (9.95, 26.667],
      (9.95, 26.667], (9.95, 26.667], (26.667, 43.333], (43.333, 60],
      (43.333, 60]], dtype=object)
```

```
Levels (3): Index([(9.95, 26.667], (26.667, 43.333], (43.333, 60]], dtype=object)
```

If the `bins` keyword is an integer, then equal-width bins are formed. Alternatively we can specify custom bin-edges:

```
In [1043]: cut(ages, bins=[0, 18, 35, 70])
```

```
Out[1043]:
```

```
Categorical:
```

```
array([(0, 18], (0, 18], (0, 18], (0, 18], (18, 35], (18, 35], (18, 35],
      (35, 70], (35, 70]], dtype=object)
```

```
Levels (3): Index([(0, 18], (18, 35], (35, 70]], dtype=object)
```



# TIME SERIES / DATE FUNCTIONALITY

pandas has proven very successful as a tool for working with time series data, especially in the financial data analysis space. With the 0.8 release, we have further improved the time series API in pandas by leaps and bounds. Using the new NumPy `datetime64` dtype, we have consolidated a large number of features from other Python libraries like `scikits.timeseries` as well as created a tremendous amount of new functionality for manipulating time series data.

In working with time series data, we will frequently seek to:

- generate sequences of fixed-frequency dates and time spans
- conform or convert time series to a particular frequency
- compute “relative” dates based on various non-standard time increments (e.g. 5 business days before the last business day of the year), or “roll” dates forward or backward

pandas provides a relatively compact and self-contained set of tools for performing the above tasks.

Create a range of dates:

```
# 72 hours starting with midnight Jan 1st, 2011
In [1067]: rng = date_range('1/1/2011', periods=72, freq='H')

In [1068]: rng[:5]
Out[1068]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-01-01 00:00:00, ..., 2011-01-01 04:00:00]
Length: 5, Freq: H, Timezone: None
```

Index pandas objects with dates:

```
In [1069]: ts = Series(randn(len(rng)), index=rng)

In [1070]: ts.head()
Out[1070]:
2011-01-01 00:00:00    0.469112
2011-01-01 01:00:00   -0.282863
2011-01-01 02:00:00   -1.509059
2011-01-01 03:00:00   -1.135632
2011-01-01 04:00:00    1.212112
Freq: H
```

Change frequency and fill gaps:

```
# to 45 minute frequency and forward fill
In [1071]: converted = ts.asfreq('45Min', method='pad')
```

```
In [1072]: converted.head()
Out[1072]:
2011-01-01 00:00:00    0.469112
2011-01-01 00:45:00    0.469112
2011-01-01 01:30:00   -0.282863
2011-01-01 02:15:00   -1.509059
2011-01-01 03:00:00   -1.135632
Freq: 45T
```

Resample:

```
# Daily means
In [1073]: ts.resample('D', how='mean')
Out[1073]:
2011-01-01    0.469112
2011-01-02   -0.322252
2011-01-03   -0.317244
2011-01-04    0.083412
Freq: D
```

## 13.1 Time Stamps vs. Time Spans

Time-stamped data is the most basic type of timeseries data that associates values with points in time. For pandas objects it means using the points in time to create the index

```
In [1074]: dates = [datetime(2012, 5, 1), datetime(2012, 5, 2), datetime(2012, 5, 3)]

In [1075]: ts = Series(np.random.randn(3), dates)

In [1076]: type(ts.index)
Out[1076]: pandas.tseries.index.DatetimeIndex

In [1077]: ts
Out[1077]:
2012-05-01   -0.410001
2012-05-02   -0.078638
2012-05-03    0.545952
```

However, in many cases it is more natural to associate things like change variables with a time span instead.

For example:

```
In [1078]: periods = PeriodIndex([Period('2012-01'), Period('2012-02'),
.....:                           Period('2012-03')])
.....:

In [1079]: ts = Series(np.random.randn(3), periods)

In [1080]: type(ts.index)
Out[1080]: pandas.tseries.period.PeriodIndex

In [1081]: ts
Out[1081]:
Jan-2012   -1.219217
Feb-2012   -1.226825
Mar-2012    0.769804
Freq: M
```



Starting with 0.8, pandas allows you to capture both representations and convert between them. Under the hood, pandas represents timestamps using instances of `Timestamp` and sequences of timestamps using instances of `DatetimeIndex`. For regular time spans, pandas uses `Period` objects for scalar values and `PeriodIndex` for sequences of spans. Better support for irregular intervals with arbitrary start and end points are forth-coming in future releases.

## 13.2 Generating Ranges of Timestamps

To generate an index with time stamps, you can use either the `DatetimeIndex` or `Index` constructor and pass in a list of datetime objects:

```
In [1082]: dates = [datetime(2012, 5, 1), datetime(2012, 5, 2), datetime(2012, 5, 3)]
```

```
In [1083]: index = DatetimeIndex(dates)
```

```
In [1084]: index # Note the frequency information
```

```
Out[1084]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2012-05-01 00:00:00, ..., 2012-05-03 00:00:00]
Length: 3, Freq: None, Timezone: None
```

```
In [1085]: index = Index(dates)
```

```
In [1086]: index # Automatically converted to DatetimeIndex
```

```
Out[1086]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2012-05-01 00:00:00, ..., 2012-05-03 00:00:00]
Length: 3, Freq: None, Timezone: None
```

Practically, this becomes very cumbersome because we often need a very long index with a large number of timestamps. If we need timestamps on a regular frequency, we can use the pandas functions `date_range` and `bdate_range` to create timestamp indexes.

```
In [1087]: index = date_range('2000-1-1', periods=1000, freq='M')
```

```
In [1088]: index
```

```
Out[1088]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2000-01-31 00:00:00, ..., 2083-04-30 00:00:00]
Length: 1000, Freq: M, Timezone: None
```

```
In [1089]: index = bdate_range('2012-1-1', periods=250)
```

```
In [1090]: index
```

```
Out[1090]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2012-01-02 00:00:00, ..., 2012-12-14 00:00:00]
Length: 250, Freq: B, Timezone: None
```

Convenience functions like `date_range` and `bdate_range` utilizes a variety of frequency aliases. The default frequency for `date_range` is a **calendar day** while the default for `bdate_range` is a **business day**

```
In [1091]: start = datetime(2011, 1, 1)
```

```
In [1092]: end = datetime(2012, 1, 1)
```

```
In [1093]: rng = date_range(start, end)
```

```
In [1094]: rng
Out[1094]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-01-01 00:00:00, ..., 2012-01-01 00:00:00]
Length: 366, Freq: D, Timezone: None
```

```
In [1095]: rng = bdate_range(start, end)
```

```
In [1096]: rng
Out[1096]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-01-03 00:00:00, ..., 2011-12-30 00:00:00]
Length: 260, Freq: B, Timezone: None
```

`date_range` and `bdate_range` makes it easy to generate a range of dates using various combinations of its parameters like `start`, `end`, `periods`, and `freq`:

```
In [1097]: date_range(start, end, freq='BM')
Out[1097]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-01-31 00:00:00, ..., 2011-12-30 00:00:00]
Length: 12, Freq: BM, Timezone: None
```

```
In [1098]: date_range(start, end, freq='W')
Out[1098]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-01-02 00:00:00, ..., 2012-01-01 00:00:00]
Length: 53, Freq: W-SUN, Timezone: None
```

```
In [1099]: bdate_range(end=end, periods=20)
Out[1099]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-12-05 00:00:00, ..., 2011-12-30 00:00:00]
Length: 20, Freq: B, Timezone: None
```

```
In [1100]: bdate_range(start=start, periods=20)
Out[1100]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-01-03 00:00:00, ..., 2011-01-28 00:00:00]
Length: 20, Freq: B, Timezone: None
```

The start and end dates are strictly inclusive. So it will not generate any dates outside of those dates if specified.

### 13.2.1 DatetimeIndex

One of the main uses for `DatetimeIndex` is as an index for pandas objects. The `DatetimeIndex` class contains many timeseries related optimizations:

- A large range of dates for various offsets are pre-computed and cached under the hood in order to make generating subsequent date ranges very fast (just have to grab a slice)
- Fast shifting using the `shift` and `tshift` method on pandas objects
- Unioning of overlapping `DatetimeIndex` objects with the same frequency is very fast (important for fast data alignment)
- Quick access to date fields via properties such as `year`, `month`, etc.

- Regularization functions like `snap` and very fast `asof` logic

`DatetimeIndex` can be used like a regular index and offers all of its intelligent functionality like selection, slicing, etc.

```
In [1101]: rng = date_range(start, end, freq='BM')
```

```
In [1102]: ts = Series(randn(len(rng)), index=rng)
```

```
In [1103]: ts.index
```

```
Out[1103]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-01-31 00:00:00, ..., 2011-12-30 00:00:00]
Length: 12, Freq: BM, Timezone: None
```

```
In [1104]: ts[:5].index
```

```
Out[1104]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-01-31 00:00:00, ..., 2011-05-31 00:00:00]
Length: 5, Freq: BM, Timezone: None
```

```
In [1105]: ts[::2].index
```

```
Out[1105]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-01-31 00:00:00, ..., 2011-11-30 00:00:00]
Length: 6, Freq: 2BM, Timezone: None
```

You can pass in dates and strings that parses to dates as indexing parameters:

```
In [1106]: ts['1/31/2011']
```

```
Out[1106]: -1.2812473076599531
```

```
In [1107]: ts[datetime(2011, 12, 25):]
```

```
Out[1107]:
2011-12-30    0.687738
Freq: BM
```

```
In [1108]: ts['10/31/2011':'12/31/2011']
```

```
Out[1108]:
2011-10-31    0.149748
2011-11-30   -0.732339
2011-12-30    0.687738
Freq: BM
```

A truncate convenience function is provided that is equivalent to slicing:

```
In [1109]: ts.truncate(before='10/31/2011', after='12/31/2011')
```

```
Out[1109]:
2011-10-31    0.149748
2011-11-30   -0.732339
2011-12-30    0.687738
Freq: BM
```

To provide convenience for accessing longer time series, you can also pass in the year or year and month as strings:

```
In [1110]: ts['2011']
```

```
Out[1110]:
2011-01-31   -1.281247
2011-02-28   -0.727707
2011-03-31   -0.121306
```

```
2011-04-29    -0.097883
2011-05-31     0.695775
2011-06-30     0.341734
2011-07-29     0.959726
2011-08-31    -1.110336
2011-09-30    -0.619976
2011-10-31     0.149748
2011-11-30    -0.732339
2011-12-30     0.687738
Freq: BM
```

```
In [1111]: ts['2011-6']
Out[1111]:
2011-06-30     0.341734
Freq: BM
```

Even complicated fancy indexing that breaks the `DatetimeIndex`'s frequency regularity will result in a `DatetimeIndex` (but frequency is lost):

```
In [1112]: ts[[0, 2, 6]].index
Out[1112]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-01-31 00:00:00, ..., 2011-07-29 00:00:00]
Length: 3, Freq: None, Timezone: None
```

`DatetimeIndex` objects has all the basic functionality of regular `Index` objects and a smorgasbord of advanced timeseries-specific methods for easy frequency processing.

**See Also:**

*[Reindexing methods](#)*

---

**Note:** While pandas does not force you to have a sorted date index, some of these methods may have unexpected or incorrect behavior if the dates are unsorted. So please be careful.

---

## 13.3 DateOffset objects

In the preceding examples, we created `DatetimeIndex` objects at various frequencies by passing in frequency strings like 'M', 'W', and 'BM' to the `freq` keyword. Under the hood, these frequency strings are being translated into an instance of pandas `DateOffset`, which represents a regular frequency increment. Specific offset logic like "month", "business day", or "one hour" is represented in its various subclasses.

Class name	Description
DateOffset	Generic offset class, defaults to 1 calendar day
BDay	business day (weekday)
Week	one week, optionally anchored on a day of the week
WeekOfMonth	the x-th day of the y-th week of each month
MonthEnd	calendar month end
MonthBegin	calendar month begin
BMonthEnd	business month end
BMonthBegin	business month begin
QuarterEnd	calendar quarter end
QuarterBegin	calendar quarter begin
BQuarterEnd	business quarter end
BQuarterBegin	business quarter begin
YearEnd	calendar year end
YearBegin	calendar year begin
BYearEnd	business year end
BYearBegin	business year begin
Hour	one hour
Minute	one minute
Second	one second
Milli	one millisecond
Micro	one microsecond

The basic `DateOffset` takes the same arguments as `datetime.timedelta`, which works like:

```
In [1113]: d = datetime(2008, 8, 18)

In [1114]: d + relativedelta(months=4, days=5)
Out[1114]: datetime.datetime(2008, 12, 23, 0, 0)
```

We could have done the same thing with `DateOffset`:

```
In [1115]: from pandas.tseries.offsets import *

In [1116]: d + DateOffset(months=4, days=5)
Out[1116]: datetime.datetime(2008, 12, 23, 0, 0)
```

The key features of a `DateOffset` object are:

- it can be added / subtracted to/from a datetime object to obtain a shifted date
- it can be multiplied by an integer (positive or negative) so that the increment will be applied multiple times
- it has `rollforward` and `rollback` methods for moving a date forward or backward to the next or previous “offset date”

Subclasses of `DateOffset` define the `apply` function which dictates custom date increment logic, such as adding business days:

```
class BDay(DateOffset):
    """DateOffset increments between business days"""
    def apply(self, other):
        ...

In [1117]: d - 5 * BDay()
Out[1117]: datetime.datetime(2008, 8, 11, 0, 0)

In [1118]: d + BMonthEnd()
Out[1118]: datetime.datetime(2008, 8, 29, 0, 0)
```

The `rollforward` and `rollback` methods do exactly what you would expect:

```
In [1119]: d
Out[1119]: datetime.datetime(2008, 8, 18, 0, 0)

In [1120]: offset = BMonthEnd()

In [1121]: offset.rollforward(d)
Out[1121]: datetime.datetime(2008, 8, 29, 0, 0)

In [1122]: offset.rollback(d)
Out[1122]: datetime.datetime(2008, 7, 31, 0, 0)
```

It's definitely worth exploring the `pandas.tseries.offsets` module and the various docstrings for the classes.

### 13.3.1 Parametric offsets

Some of the offsets can be “parameterized” when created to result in different behavior. For example, the `Week` offset for generating weekly data accepts a `weekday` parameter which results in the generated dates always lying on a particular day of the week:

```
In [1123]: d + Week()
Out[1123]: datetime.datetime(2008, 8, 25, 0, 0)

In [1124]: d + Week(weekday=4)
Out[1124]: datetime.datetime(2008, 8, 22, 0, 0)

In [1125]: (d + Week(weekday=4)).weekday()
Out[1125]: 4
```

Another example is parameterizing `YearEnd` with the specific ending month:

```
In [1126]: d + YearEnd()
Out[1126]: datetime.datetime(2008, 12, 31, 0, 0)

In [1127]: d + YearEnd(month=6)
Out[1127]: datetime.datetime(2009, 6, 30, 0, 0)
```

### 13.3.2 Offset Aliases

A number of string aliases are given to useful common time series frequencies. We will refer to these aliases as *offset aliases* (referred to as *time rules* prior to v0.8.0).

Alias	Description
B	business day frequency
D	calendar day frequency
W	weekly frequency
M	month end frequency
BM	business month end frequency
MS	month start frequency
BMS	business month start frequency
Q	quarter end frequency
BQ	business quarter end frequency
QS	quarter start frequency
BQS	business quarter start frequency
A	year end frequency
BA	business year end frequency
AS	year start frequency
BAS	business year start frequency
H	hourly frequency
T	minutely frequency
S	secondly frequency
L	milliseconds
U	microseconds

### 13.3.3 Combining Aliases

As we have seen previously, the alias and the offset instance are fungible in most functions:

```
In [1128]: date_range(start, periods=5, freq='B')
Out[1128]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-01-03 00:00:00, ..., 2011-01-07 00:00:00]
Length: 5, Freq: B, Timezone: None
```

```
In [1129]: date_range(start, periods=5, freq=BDay())
Out[1129]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-01-03 00:00:00, ..., 2011-01-07 00:00:00]
Length: 5, Freq: B, Timezone: None
```

You can combine together day and intraday offsets:

```
In [1130]: date_range(start, periods=10, freq='2h20min')
Out[1130]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-01-01 00:00:00, ..., 2011-01-01 21:00:00]
Length: 10, Freq: 140T, Timezone: None
```

```
In [1131]: date_range(start, periods=10, freq='1D10U')
Out[1131]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2011-01-01 00:00:00, ..., 2011-01-10 00:00:00.000090]
Length: 10, Freq: 86400000010U, Timezone: None
```

### 13.3.4 Anchored Offsets

For some frequencies you can specify an anchoring suffix:

Alias	Description
W-SUN	weekly frequency (sundays). Same as 'W'
W-MON	weekly frequency (mondays)
W-TUE	weekly frequency (tuesdays)
W-WED	weekly frequency (wednesdays)
W-THU	weekly frequency (thursdays)
W-FRI	weekly frequency (fridays)
W-SAT	weekly frequency (saturdays)
(B)Q(S)-DEC	quarterly frequency, year ends in December. Same as 'Q'
(B)Q(S)-JAN	quarterly frequency, year ends in January
(B)Q(S)-FEB	quarterly frequency, year ends in February
(B)Q(S)-MAR	quarterly frequency, year ends in March
(B)Q(S)-APR	quarterly frequency, year ends in April
(B)Q(S)-MAY	quarterly frequency, year ends in May
(B)Q(S)-JUN	quarterly frequency, year ends in June
(B)Q(S)-JUL	quarterly frequency, year ends in July
(B)Q(S)-AUG	quarterly frequency, year ends in August
(B)Q(S)-SEP	quarterly frequency, year ends in September
(B)Q(S)-OCT	quarterly frequency, year ends in October
(B)Q(S)-NOV	quarterly frequency, year ends in November
(B)A(S)-DEC	annual frequency, anchored end of December. Same as 'A'
(B)A(S)-JAN	annual frequency, anchored end of January
(B)A(S)-FEB	annual frequency, anchored end of February
(B)A(S)-MAR	annual frequency, anchored end of March
(B)A(S)-APR	annual frequency, anchored end of April
(B)A(S)-MAY	annual frequency, anchored end of May
(B)A(S)-JUN	annual frequency, anchored end of June
(B)A(S)-JUL	annual frequency, anchored end of July
(B)A(S)-AUG	annual frequency, anchored end of August
(B)A(S)-SEP	annual frequency, anchored end of September
(B)A(S)-OCT	annual frequency, anchored end of October
(B)A(S)-NOV	annual frequency, anchored end of November

These can be used as arguments to `date_range`, `bdate_range`, constructors for `DatetimeIndex`, as well as various other timeseries-related functions in pandas.

### 13.3.5 Legacy Aliases

Note that prior to v0.8.0, time rules had a slightly different look. Pandas will continue to support the legacy time rules for the time being but it is strongly recommended that you switch to using the new offset aliases.



Legacy Time Rule	Offset Alias
WEEKDAY	B
EOM	BM
W@MON	W-MON
W@TUE	W-TUE
W@WED	W-WED
W@THU	W-THU
W@FRI	W-FRI
W@SAT	W-SAT
W@SUN	W-SUN
Q@JAN	BQ-JAN
Q@FEB	BQ-FEB
Q@MAR	BQ-MAR
A@JAN	BA-JAN
A@FEB	BA-FEB
A@MAR	BA-MAR
A@APR	BA-APR
A@MAY	BA-MAY
A@JUN	BA-JUN
A@JUL	BA-JUL
A@AUG	BA-AUG
A@SEP	BA-SEP
A@OCT	BA-OCT
A@NOV	BA-NOV
A@DEC	BA-DEC
min	T
ms	L
us: “U”	

As you can see, legacy quarterly and annual frequencies are business quarter and business year ends. Please also note the legacy time rule for milliseconds `ms` versus the new offset alias for month start `MS`. This means that offset alias parsing is case sensitive.

## 13.4 Time series-related instance methods

### 13.4.1 Shifting / lagging

One may want to *shift* or *lag* the values in a `TimeSeries` back and forward in time. The method for this is `shift`, which is available on all of the pandas objects. In `DataFrame`, `shift` will currently only shift along the `index` and in `Panel` along the `major_axis`.

```
In [1132]: ts = ts[:5]
```

```
In [1133]: ts.shift(1)
```

```
Out[1133]:
2011-01-31      NaN
2011-02-28    -1.281247
2011-03-31    -0.727707
2011-04-29    -0.121306
2011-05-31    -0.097883
Freq: BM
```

The `shift` method accepts an `freq` argument which can accept a `DateOffset` class or other `timedelta`-like object

or also a *offset alias*:

```
In [1134]: ts.shift(5, freq=datetools.bday)
Out[1134]:
2011-02-07    -1.281247
2011-03-07    -0.727707
2011-04-07    -0.121306
2011-05-06    -0.097883
2011-06-07     0.695775

In [1135]: ts.shift(5, freq='BM')
Out[1135]:
2011-06-30    -1.281247
2011-07-29    -0.727707
2011-08-31    -0.121306
2011-09-30    -0.097883
2011-10-31     0.695775
Freq: BM
```

Rather than changing the alignment of the data and the index, `DataFrame` and `TimeSeries` objects also have a `tshift` convenience method that changes all the dates in the index by a specified number of offsets:

```
In [1136]: ts.tshift(5, freq='D')
Out[1136]:
2011-02-05    -1.281247
2011-03-05    -0.727707
2011-04-05    -0.121306
2011-05-04    -0.097883
2011-06-05     0.695775
```

Note that with `tshift`, the leading entry is no longer `NaN` because the data is not being realigned.

## 13.4.2 Frequency conversion

The primary function for changing frequencies is the `asfreq` function. For a `DatetimeIndex`, this is basically just a thin, but convenient wrapper around `reindex` which generates a `date_range` and calls `reindex`.

```
In [1137]: dr = date_range('1/1/2010', periods=3, freq=3 * datetools.bday)
```

```
In [1138]: ts = Series(randn(3), index=dr)
```

```
In [1139]: ts
Out[1139]:
2010-01-01    0.176444
2010-01-06    0.403310
2010-01-11   -0.154951
Freq: 3B
```

```
In [1140]: ts.asfreq(BDay())
Out[1140]:
2010-01-01    0.176444
2010-01-04         NaN
2010-01-05         NaN
2010-01-06    0.403310
2010-01-07         NaN
2010-01-08         NaN
2010-01-11   -0.154951
Freq: B
```

`asfreq` provides a further convenience so you can specify an interpolation method for any gaps that may appear after the frequency conversion

```
In [1141]: ts.asfreq(BDay(), method='pad')
Out[1141]:
2010-01-01    0.176444
2010-01-04    0.176444
2010-01-05    0.176444
2010-01-06    0.403310
2010-01-07    0.403310
2010-01-08    0.403310
2010-01-11   -0.154951
Freq: B
```

### 13.4.3 Filling forward / backward

Related to `asfreq` and `reindex` is the `fillna` function documented in the [missing data section](#).

## 13.5 Up- and downsampling

With 0.8, pandas introduces simple, powerful, and efficient functionality for performing resampling operations during frequency conversion (e.g., converting secondly data into 5-minutely data). This is extremely common in, but not limited to, financial applications.

```
In [1142]: rng = date_range('1/1/2012', periods=100, freq='S')
In [1143]: ts = Series(randint(0, 500, len(rng)), index=rng)

In [1144]: ts.resample('5Min', how='sum')
Out[1144]:
2012-01-01 00:00:00    230
2012-01-01 00:05:00   25562
Freq: 5T
```

The `resample` function is very flexible and allows you to specify many different parameters to control the frequency conversion and resampling operation.

The `how` parameter can be a function name or numpy array function that takes an array and produces an aggregated value:

```
In [1145]: ts.resample('5Min') # default is mean
Out[1145]:
2012-01-01 00:00:00    230.00000
2012-01-01 00:05:00   258.20202
Freq: 5T

In [1146]: ts.resample('5Min', how='ohlc')
Out[1146]:
              open  high  low  close
2012-01-01 00:00:00   230   230   230    230
2012-01-01 00:05:00   202   492    0    214

In [1147]: ts.resample('5Min', how=np.max)
Out[1147]:
2012-01-01 00:00:00    230
2012-01-01 00:05:00   492
```

Any function available via *dispatching* can be given to the `how` parameter by name, including `sum`, `mean`, `std`, `max`, `min`, `median`, `first`, `last`, `ohlc`.

For downsampling, `closed` can be set to `'left'` or `'right'` to specify which end of the interval is closed:

```
In [1148]: ts.resample('5Min', closed='right')
Out[1148]:
2012-01-01 00:00:00    230.00000
2012-01-01 00:05:00    258.20202
Freq: 5T
```

```
In [1149]: ts.resample('5Min', closed='left')
Out[1149]:
2012-01-01 00:05:00    257.92
Freq: 5T
```

For upsampling, the `fill_method` and `limit` parameters can be specified to interpolate over the gaps that are created:

```
# from secondly to every 250 milliseconds
```

```
In [1150]: ts[:2].resample('250L')
Out[1150]:
2012-01-01 00:00:00    230
2012-01-01 00:00:00.250000    NaN
2012-01-01 00:00:00.500000    NaN
2012-01-01 00:00:00.750000    NaN
2012-01-01 00:00:01    202
Freq: 250L
```

```
In [1151]: ts[:2].resample('250L', fill_method='pad')
Out[1151]:
2012-01-01 00:00:00    230
2012-01-01 00:00:00.250000    230
2012-01-01 00:00:00.500000    230
2012-01-01 00:00:00.750000    230
2012-01-01 00:00:01    202
Freq: 250L
```

```
In [1152]: ts[:2].resample('250L', fill_method='pad', limit=2)
Out[1152]:
2012-01-01 00:00:00    230
2012-01-01 00:00:00.250000    230
2012-01-01 00:00:00.500000    230
2012-01-01 00:00:00.750000    NaN
2012-01-01 00:00:01    202
Freq: 250L
```

Parameters like `label` and `loffset` are used to manipulate the resulting labels. `label` specifies whether the result is labeled with the beginning or the end of the interval. `loffset` performs a time adjustment on the output labels.

```
In [1153]: ts.resample('5Min') # by default label='right'
Out[1153]:
2012-01-01 00:00:00    230.00000
2012-01-01 00:05:00    258.20202
Freq: 5T
```

```
In [1154]: ts.resample('5Min', label='left')
Out[1154]:
2011-12-31 23:55:00    230.00000
2012-01-01 00:00:00    258.20202
```

```
Freq: 5T
```

```
In [1155]: ts.resample('5Min', label='left', loffset='1s')
```

```
Out[1155]:
```

```
2011-12-31 23:55:01    230.00000
```

```
2012-01-01 00:00:01    258.20202
```

The `axis` parameter can be set to 0 or 1 and allows you to resample the specified axis for a `DataFrame`.

`kind` can be set to 'timestamp' or 'period' to convert the resulting index to/from time-stamp and time-span representations. By default `resample` retains the input representation.

`convention` can be set to 'start' or 'end' when resampling period data (detail below). It specifies how low frequency periods are converted to higher frequency periods.

Note that 0.8 marks a watershed in the timeseries functionality in pandas. In previous versions, resampling had to be done using a combination of `date_range`, `groupby` with `asof`, and then calling an aggregation function on the grouped object. This was not nearly convenient or performant as the new pandas timeseries API.

## 13.6 Time Span Representation

Regular intervals of time are represented by `Period` objects in pandas while sequences of `Period` objects are collected in a `PeriodIndex`, which can be created with the convenience function `period_range`.

### 13.6.1 Period

A `Period` represents a span of time (e.g., a day, a month, a quarter, etc). It can be created using a frequency alias:

```
In [1156]: Period('2012', freq='A-DEC')
```

```
Out[1156]: Period('2012', 'A-DEC')
```

```
In [1157]: Period('2012-1-1', freq='D')
```

```
Out[1157]: Period('01-Jan-2012', 'D')
```

```
In [1158]: Period('2012-1-1 19:00', freq='H')
```

```
Out[1158]: Period('01-Jan-2012 19:00', 'H')
```

Unlike time stamped data, pandas does not support frequencies at multiples of `DateOffsets` (e.g., '3Min') for periods.

Adding and subtracting integers from periods shifts the period by its own frequency.

```
In [1159]: p = Period('2012', freq='A-DEC')
```

```
In [1160]: p + 1
```

```
Out[1160]: Period('2013', 'A-DEC')
```

```
In [1161]: p - 3
```

```
Out[1161]: Period('2009', 'A-DEC')
```

Taking the difference of `Period` instances with the same frequency will return the number of frequency units between them:

```
In [1162]: Period('2012', freq='A-DEC') - Period(2002', freq='A-DEC')
```

```
File "<ipython-input-1162-41a08553f136>", line 1
```

```
Period('2012', freq='A-DEC') - Period(2002', freq='A-DEC')
```

```
^
```

```
SyntaxError: invalid syntax
```

### 13.6.2 PeriodIndex and period\_range

Regular sequences of Period objects can be collected in a PeriodIndex, which can be constructed using the period\_range convenience function:

```
In [1163]: prng = period_range('1/1/2011', '1/1/2012', freq='M')
```

```
In [1164]: prng
Out[1164]:
<class 'pandas.tseries.period.PeriodIndex'>
freq: M
[Jan-2011, ..., Jan-2012]
length: 13
```

The PeriodIndex constructor can also be used directly:

```
In [1165]: PeriodIndex(['2011-1', '2011-2', '2011-3'], freq='M')
Out[1165]:
<class 'pandas.tseries.period.PeriodIndex'>
freq: M
[Jan-2011, ..., Mar-2011]
length: 3
```

Just like DatetimeIndex, a PeriodIndex can also be used to index pandas objects:

```
In [1166]: Series(randn(len(prng)), prng)
Out[1166]:
Jan-2011    0.301624
Feb-2011   -1.460489
Mar-2011    0.610679
Apr-2011    1.195856
May-2011   -0.008820
Jun-2011   -0.045729
Jul-2011   -1.051015
Aug-2011   -0.422924
Sep-2011   -0.028361
Oct-2011   -0.782386
Nov-2011    0.861980
Dec-2011    1.438604
Jan-2012   -0.525492
Freq: M
```

### 13.6.3 Frequency Conversion and Resampling with PeriodIndex

The frequency of Periods and PeriodIndex can be converted via the asfreq method. Let's start with the fiscal year 2011, ending in December:

```
In [1167]: p = Period('2011', freq='A-DEC')
```

```
In [1168]: p
Out[1168]: Period('2011', 'A-DEC')
```

We can convert it to a monthly frequency. Using the how parameter, we can specify whether to return the starting or ending month:

```
In [1169]: p.asfreq('M', how='start')
Out[1169]: Period('Jan-2011', 'M')
```

```
In [1170]: p.asfreq('M', how='end')
Out[1170]: Period('Dec-2011', 'M')
```

The shorthands 's' and 'e' are provided for convenience:

```
In [1171]: p.asfreq('M', 's')
Out[1171]: Period('Jan-2011', 'M')
```

```
In [1172]: p.asfreq('M', 'e')
Out[1172]: Period('Dec-2011', 'M')
```

Converting to a “super-period” (e.g., annual frequency is a super-period of quarterly frequency) automatically returns the super-period that includes the input period:

```
In [1173]: p = Period('2011-12', freq='M')
```

```
In [1174]: p.asfreq('A-NOV')
Out[1174]: Period('2012', 'A-NOV')
```

Note that since we converted to an annual frequency that ends the year in November, the monthly period of December 2011 is actually in the 2012 A-NOV period. Period conversions with anchored frequencies are particularly useful for working with various quarterly data common to economics, business, and other fields. Many organizations define quarters relative to the month in which their fiscal year start and ends. Thus, first quarter of 2011 could start in 2010 or a few months into 2011. Via anchored frequencies, pandas works all quarterly frequencies Q-JAN through Q-DEC.

Q-DEC define regular calendar quarters:

```
In [1175]: p = Period('2012Q1', freq='Q-DEC')
```

```
In [1176]: p.asfreq('D', 's')
Out[1176]: Period('01-Jan-2012', 'D')
```

```
In [1177]: p.asfreq('D', 'e')
Out[1177]: Period('31-Mar-2012', 'D')
```

Q-MAR defines fiscal year end in March:

```
In [1178]: p = Period('2011Q4', freq='Q-MAR')
```

```
In [1179]: p.asfreq('D', 's')
Out[1179]: Period('01-Jan-2011', 'D')
```

```
In [1180]: p.asfreq('D', 'e')
Out[1180]: Period('31-Mar-2011', 'D')
```

## 13.7 Converting between Representations

Timestamped data can be converted to PeriodIndex-ed data using `to_period` and vice-versa using `to_timestamp`:

```
In [1181]: rng = date_range('1/1/2012', periods=5, freq='M')
```

```
In [1182]: ts = Series(randn(len(rng)), index=rng)
```

```
In [1183]: ts
Out[1183]:
2012-01-31    -1.684469
```

```
2012-02-29      0.550605
2012-03-31      0.091955
2012-04-30      0.891713
2012-05-31      0.807078
Freq: M
```

```
In [1184]: ps = ts.to_period()
```

```
In [1185]: ps
```

```
Out [1185]:
Jan-2012    -1.684469
Feb-2012      0.550605
Mar-2012      0.091955
Apr-2012      0.891713
May-2012      0.807078
Freq: M
```

```
In [1186]: ps.to_timestamp()
```

```
Out [1186]:
2012-01-31    -1.684469
2012-02-29      0.550605
2012-03-31      0.091955
2012-04-30      0.891713
2012-05-31      0.807078
Freq: M
```

Remember that 's' and 'e' can be used to return the timestamps at the start or end of the period:

```
In [1187]: ps.to_timestamp('D', how='s')
```

```
Out [1187]:
2012-01-01    -1.684469
2012-02-01      0.550605
2012-03-01      0.091955
2012-04-01      0.891713
2012-05-01      0.807078
Freq: MS
```

Converting between period and timestamp enables some convenient arithmetic functions to be used. In the following example, we convert a quarterly frequency with year ending in November to 9am of the end of the month following the quarter end:

```
In [1188]: prng = period_range('1990Q1', '2000Q4', freq='Q-NOV')
```

```
In [1189]: ts = Series(randn(len(prng)), prng)
```

```
In [1190]: ts.index = (prng.asfreq('M', 'e') + 1).asfreq('H', 's') + 9
```

```
In [1191]: ts.head()
```

```
Out [1191]:
01-Mar-1990 09:00      0.221441
01-Jun-1990 09:00     -0.113139
01-Sep-1990 09:00     -1.812900
01-Dec-1990 09:00     -0.053708
01-Mar-1991 09:00     -0.114574
Freq: H
```



## 13.8 Time Zone Handling

Using `pytz`, pandas provides rich support for working with timestamps in different time zones. By default, pandas objects are time zone unaware:

```
In [1192]: rng = date_range('3/6/2012 00:00', periods=15, freq='D')
```

```
In [1193]: print(rng.tz)
```

```
None
```

To supply the time zone, you can use the `tz` keyword to `date_range` and other functions:

```
In [1194]: rng_utc = date_range('3/6/2012 00:00', periods=10, freq='D', tz='UTC')
```

```
In [1195]: print(rng_utc.tz)
```

```
UTC
```

Timestamps, like Python's `datetime.datetime` object can be either time zone naive or time zone aware. Naive time series and `DatetimeIndex` objects can be *localized* using `tz_localize`:

```
In [1196]: ts = Series(randn(len(rng)), rng)
```

```
In [1197]: ts_utc = ts.tz_localize('UTC')
```

```
In [1198]: ts_utc
```

```
Out[1198]:
```

```
2012-03-06 00:00:00+00:00    -0.114722
2012-03-07 00:00:00+00:00     0.168904
2012-03-08 00:00:00+00:00   -0.048048
2012-03-09 00:00:00+00:00     0.801196
2012-03-10 00:00:00+00:00     1.392071
2012-03-11 00:00:00+00:00   -0.048788
2012-03-12 00:00:00+00:00   -0.808838
2012-03-13 00:00:00+00:00   -1.003677
2012-03-14 00:00:00+00:00   -0.160766
2012-03-15 00:00:00+00:00     1.758853
2012-03-16 00:00:00+00:00     0.729195
2012-03-17 00:00:00+00:00     1.359732
2012-03-18 00:00:00+00:00     2.006296
2012-03-19 00:00:00+00:00     0.870210
2012-03-20 00:00:00+00:00     0.043464
```

```
Freq: D
```

You can use the `tz_convert` method to convert pandas objects to convert tz-aware data to another time zone:

```
In [1199]: ts_utc.tz_convert('US/Eastern')
```

```
Out[1199]:
```

```
2012-03-05 19:00:00-05:00   -0.114722
2012-03-06 19:00:00-05:00    0.168904
2012-03-07 19:00:00-05:00   -0.048048
2012-03-08 19:00:00-05:00    0.801196
2012-03-09 19:00:00-05:00    1.392071
2012-03-10 19:00:00-05:00   -0.048788
2012-03-11 20:00:00-04:00   -0.808838
2012-03-12 20:00:00-04:00   -1.003677
2012-03-13 20:00:00-04:00   -0.160766
2012-03-14 20:00:00-04:00    1.758853
2012-03-15 20:00:00-04:00    0.729195
2012-03-16 20:00:00-04:00    1.359732
```

```
2012-03-17 20:00:00-04:00    2.006296
2012-03-18 20:00:00-04:00    0.870210
2012-03-19 20:00:00-04:00    0.043464
Freq: D
```

Under the hood, all timestamps are stored in UTC. Scalar values from a `DatetimeIndex` with a time zone will have their fields (day, hour, minute) localized to the time zone. However, timestamps with the same UTC value are still considered to be equal even if they are in different time zones:

```
In [1200]: rng_eastern = rng_utc.tz_convert('US/Eastern')

In [1201]: rng_berlin = rng_utc.tz_convert('Europe/Berlin')

In [1202]: rng_eastern[5]
Out[1202]: <Timestamp: 2012-03-10 19:00:00-0500 EST, tz=US/Eastern>

In [1203]: rng_berlin[5]
Out[1203]: <Timestamp: 2012-03-11 01:00:00+0100 CET, tz=Europe/Berlin>

In [1204]: rng_eastern[5] == rng_berlin[5]
Out[1204]: True
```

Like `Series`, `DataFrame`, and `DatetimeIndex`, `Timestamps` can be converted to other time zones using `tz_convert`:

```
In [1205]: rng_eastern[5]
Out[1205]: <Timestamp: 2012-03-10 19:00:00-0500 EST, tz=US/Eastern>

In [1206]: rng_berlin[5]
Out[1206]: <Timestamp: 2012-03-11 01:00:00+0100 CET, tz=Europe/Berlin>

In [1207]: rng_eastern[5].tz_convert('Europe/Berlin')
Out[1207]: <Timestamp: 2012-03-11 01:00:00+0100 CET, tz=Europe/Berlin>
```

Localization of `Timestamps` functions just like `DatetimeIndex` and `TimeSeries`:

```
In [1208]: rng[5]
Out[1208]: <Timestamp: 2012-03-11 00:00:00>

In [1209]: rng[5].tz_localize('Asia/Shanghai')
Out[1209]: <Timestamp: 2012-03-11 00:00:00+0800 CST, tz=Asia/Shanghai>
```

Operations between `TimeSeries` in difficult time zones will yield UTC `TimeSeries`, aligning the data on the UTC timestamps:

```
In [1210]: eastern = ts_utc.tz_convert('US/Eastern')

In [1211]: berlin = ts_utc.tz_convert('Europe/Berlin')

In [1212]: result = eastern + berlin

In [1213]: result
Out[1213]:
2012-03-06 00:00:00+00:00    -0.229443
2012-03-07 00:00:00+00:00     0.337809
2012-03-08 00:00:00+00:00    -0.096096
2012-03-09 00:00:00+00:00     1.602392
2012-03-10 00:00:00+00:00     2.784142
2012-03-11 00:00:00+00:00    -0.097575
2012-03-12 00:00:00+00:00    -1.617677
```

```
2012-03-13 00:00:00+00:00    -2.007353
2012-03-14 00:00:00+00:00    -0.321532
2012-03-15 00:00:00+00:00     3.517706
2012-03-16 00:00:00+00:00     1.458389
2012-03-17 00:00:00+00:00     2.719465
2012-03-18 00:00:00+00:00     4.012592
2012-03-19 00:00:00+00:00     1.740419
2012-03-20 00:00:00+00:00     0.086928
Freq: D
```

```
In [1214]: result.index
```

```
Out[1214]:
```

```
<class 'pandas.tseries.index.DatetimeIndex'>
[2012-03-06 00:00:00, ..., 2012-03-20 00:00:00]
Length: 15, Freq: D, Timezone: UTC
```



# PLOTTING WITH MATPLOTLIB

---

**Note:** We intend to build more plotting integration with `matplotlib` as time goes on.

---

We use the standard convention for referencing the matplotlib API:

```
In [1215]: import matplotlib.pyplot as plt
```

## 14.1 Basic plotting: `plot`

The `plot` method on `Series` and `DataFrame` is just a simple wrapper around `plt.plot`:

```
In [1216]: ts = Series(randn(1000), index=date_range('1/1/2000', periods=1000))
```

```
In [1217]: ts = ts.cumsum()
```

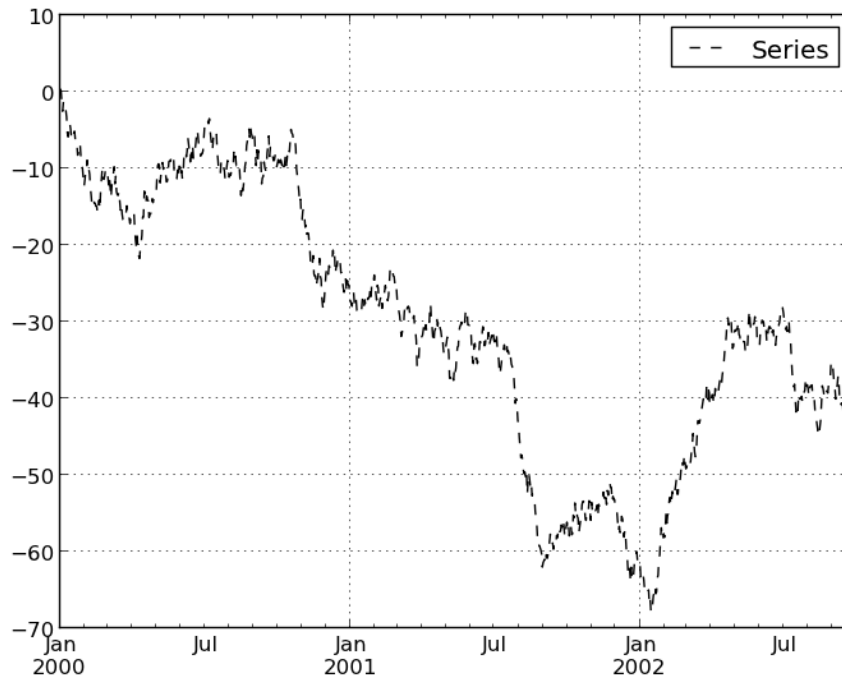
```
In [1218]: ts.plot()
```

```
Out[1218]: <matplotlib.axes.AxesSubplot at 0xba5d9d0>
```



If the index consists of dates, it calls `gcf().autofmt_xdate()` to try to format the x-axis nicely as per above. The method takes a number of arguments for controlling the look of the plot:

```
In [1219]: plt.figure(); ts.plot(style='k--', label='Series'); plt.legend()
Out[1219]: <matplotlib.legend.Legend at 0xdfa41190>
```

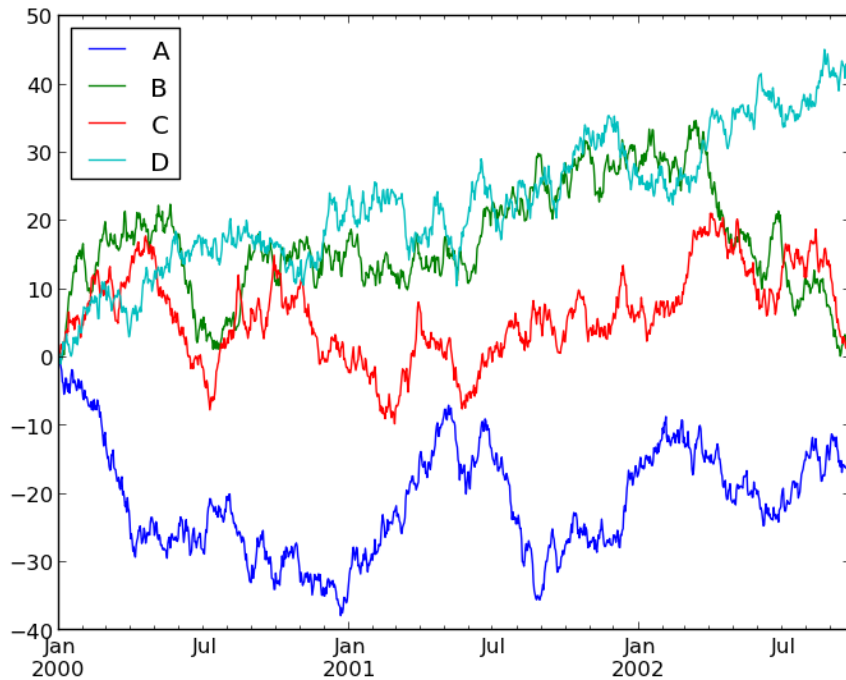


On `DataFrame`, `plot` is a convenience to plot all of the columns with labels:

```
In [1220]: df = DataFrame(randn(1000, 4), index=ts.index,
.....:                  columns=['A', 'B', 'C', 'D'])
.....:

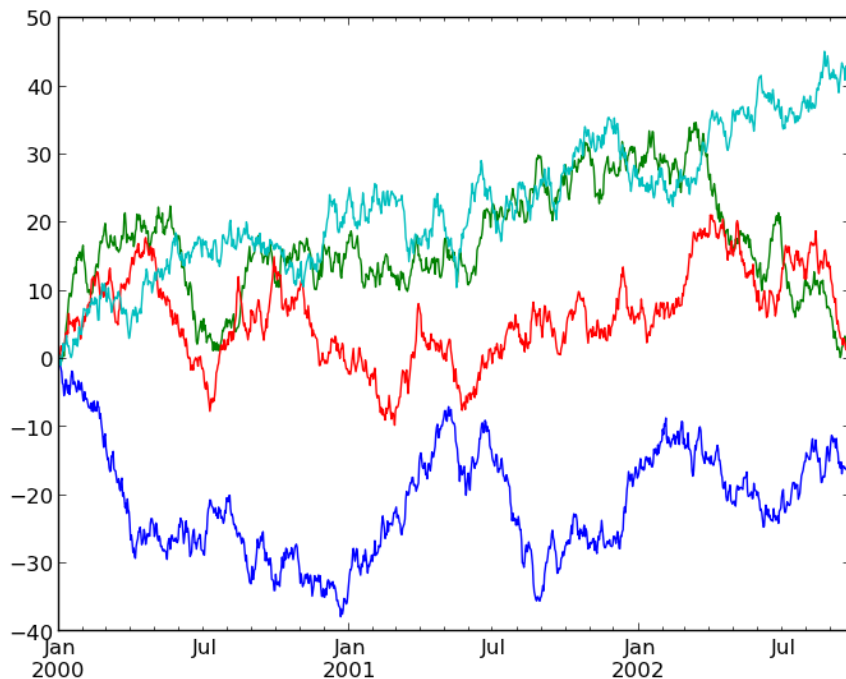
In [1221]: df = df.cumsum()

In [1222]: plt.figure(); df.plot(); plt.legend(loc='best')
Out[1222]: <matplotlib.legend.Legend at 0xdfa90d0>
```



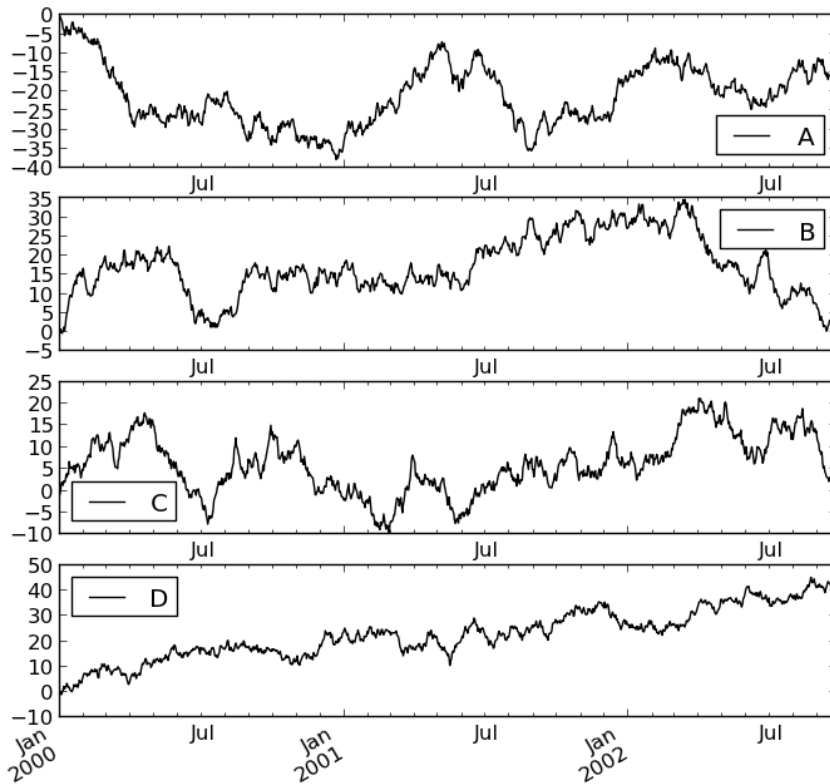
You may set the `legend` argument to `False` to hide the legend, which is shown by default.

```
In [1223]: df.plot(legend=False)
Out[1223]: <matplotlib.axes.AxesSubplot at 0xebc0c50>
```



Some other options are available, like plotting each Series on a different axis:

```
In [1224]: df.plot(subplots=True, figsize=(8, 8)); plt.legend(loc='best')
Out[1224]: <matplotlib.legend.Legend at 0xebbf690>
```



You may pass `logy` to get a log-scale Y axis.

```
In [1225]: plt.figure();
In [1225]: ts = Series(randn(1000), index=date_range('1/1/2000', periods=1000))

In [1226]: ts = np.exp(ts.cumsum())

In [1227]: ts.plot(logy=True)
Out[1227]: <matplotlib.axes.AxesSubplot at 0xfdd4d90>
```





You can plot one column versus another using the *x* and *y* keywords in *DataFrame.plot*:

```
In [1228]: plt.figure()
```

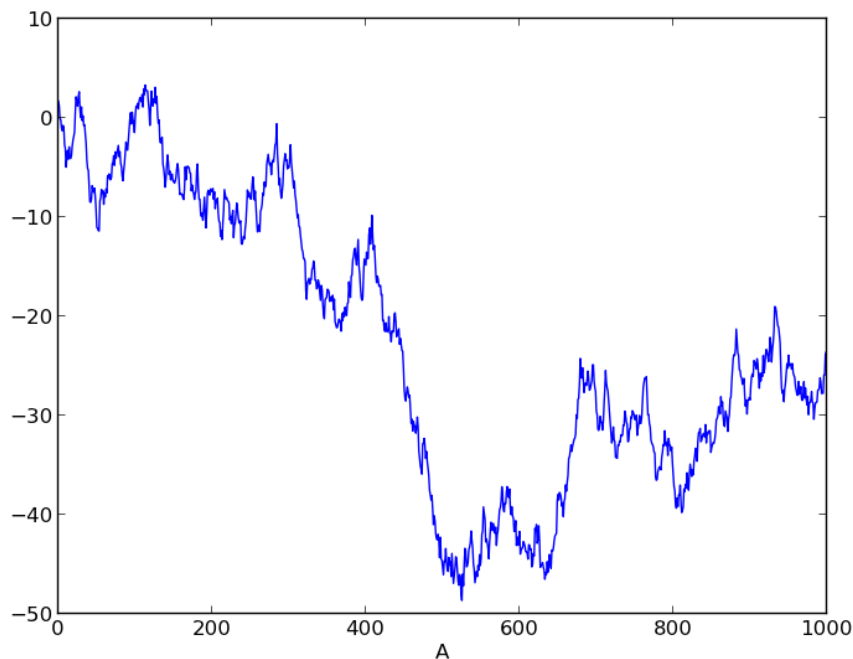
```
Out[1228]: <matplotlib.figure.Figure at 0xfdf9e90>
```

```
In [1229]: df3 = DataFrame(np.random.randn(1000, 2), columns=['B', 'C']).cumsum()
```

```
In [1230]: df3['A'] = Series(range(len(df)))
```

```
In [1231]: df3.plot(x='A', y='B')
```

```
Out[1231]: <matplotlib.axes.AxesSubplot at 0x10543050>
```



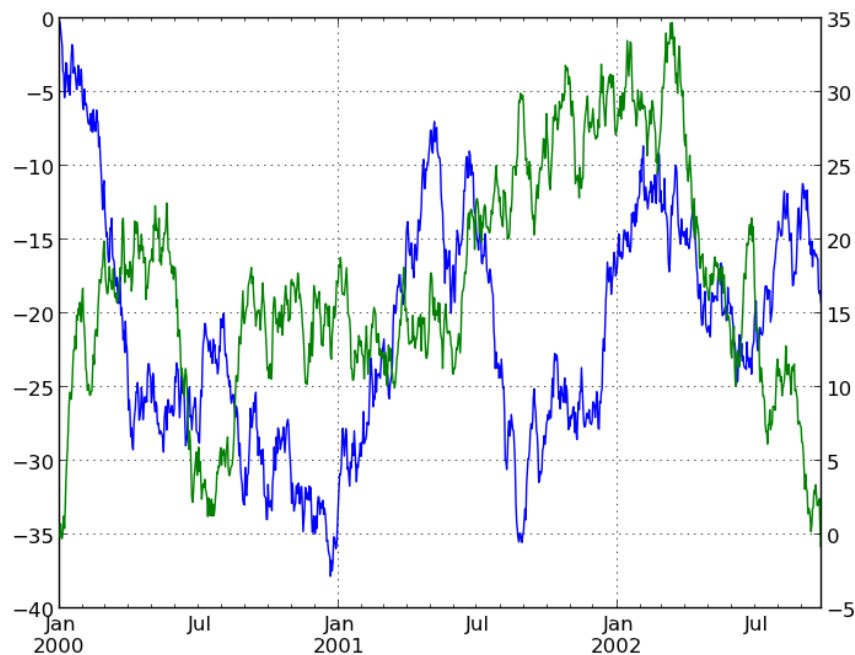
### 14.1.1 Plotting on a Secondary Y-axis

To plot data on a secondary y-axis, use the `secondary_y` keyword:

```
In [1232]: plt.figure()
Out[1232]: <matplotlib.figure.Figure at 0x1054fb50>

In [1233]: df.A.plot()
Out[1233]: <matplotlib.axes.AxesSubplot at 0x101e8250>

In [1234]: df.B.plot(secondary_y=True, style='g')
Out[1234]: <matplotlib.axes.AxesSubplot at 0x101e8250>
```

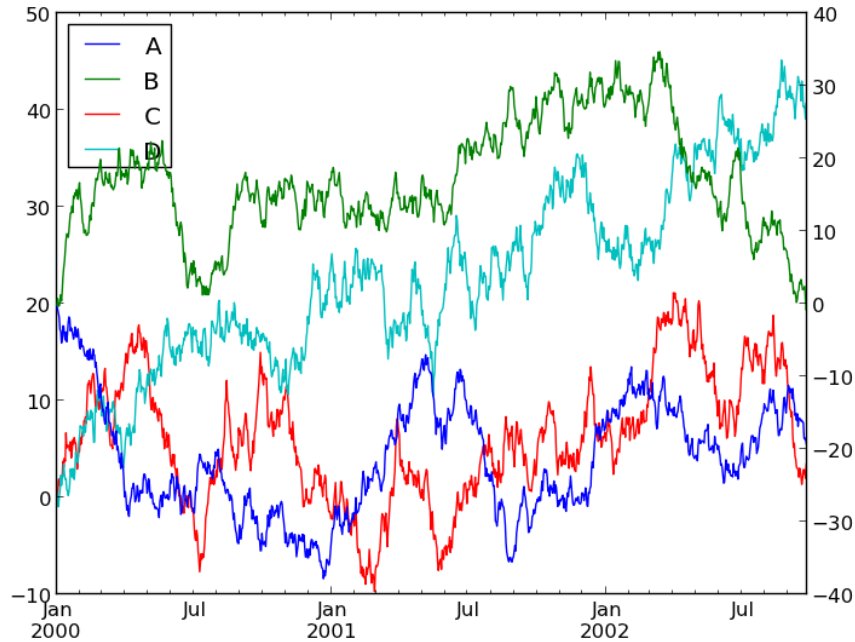


### 14.1.2 Selective Plotting on Secondary Y-axis

To plot some columns in a DataFrame, give the column names to the `secondary_y` keyword:

```
In [1235]: plt.figure()
Out[1235]: <matplotlib.figure.Figure at 0xfa27f90>

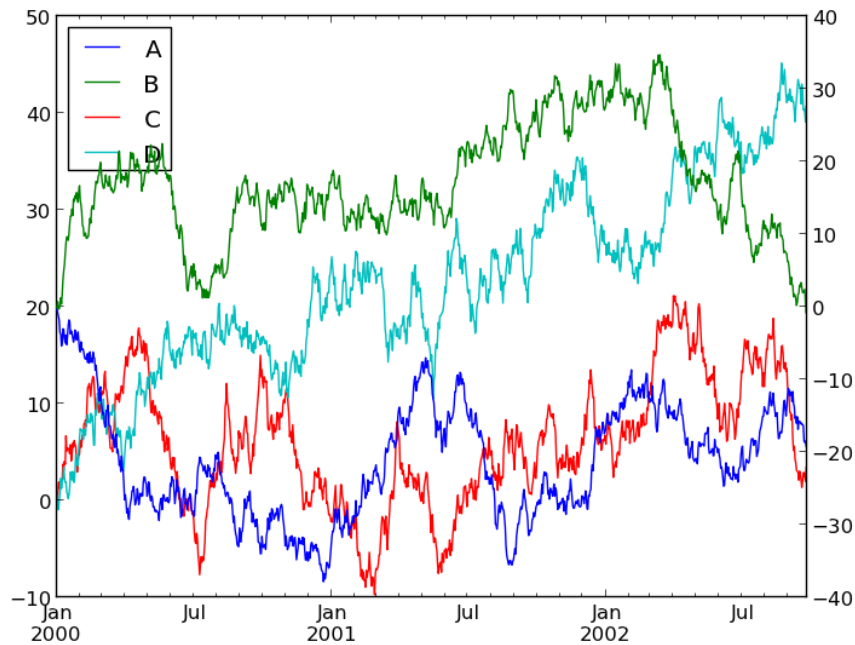
In [1236]: df.plot(secondary_y=['A', 'B'])
Out[1236]: <matplotlib.axes.AxesSubplot at 0x1099b490>
```



Note that the columns plotted on the secondary y-axis is automatically marked with “(right)” in the legend. To turn off the automatic marking, use the `mark_right=False` keyword:

```
In [1237]: plt.figure()
Out[1237]: <matplotlib.figure.Figure at 0x1099ab10>

In [1238]: df.plot(secondary_y=['A', 'B'], mark_right=False)
Out[1238]: <matplotlib.axes.AxesSubplot at 0x114412d0>
```



### 14.1.3 Targeting different subplots

You can pass an `ax` argument to `Series.plot` to plot on a particular axis:

```
In [1239]: fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(8, 5))
```

```
In [1240]: df['A'].plot(ax=axes[0,0]); axes[0,0].set_title('A')
```

```
Out[1240]: <matplotlib.text.Text at 0x11b09490>
```

```
In [1241]: df['B'].plot(ax=axes[0,1]); axes[0,1].set_title('B')
```

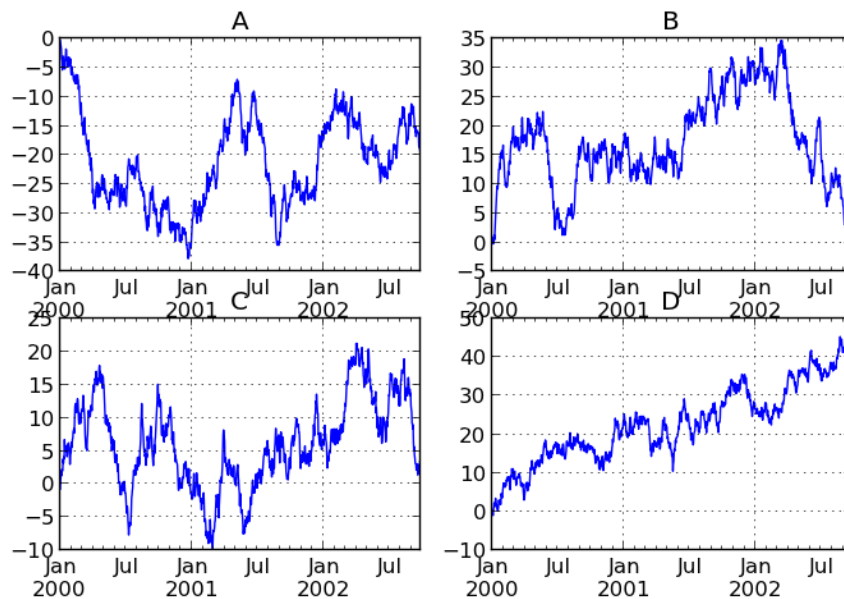
```
Out[1241]: <matplotlib.text.Text at 0x1189a090>
```

```
In [1242]: df['C'].plot(ax=axes[1,0]); axes[1,0].set_title('C')
```

```
Out[1242]: <matplotlib.text.Text at 0x11e96450>
```

```
In [1243]: df['D'].plot(ax=axes[1,1]); axes[1,1].set_title('D')
```

```
Out[1243]: <matplotlib.text.Text at 0x11ec3e10>
```



## 14.2 Other plotting features

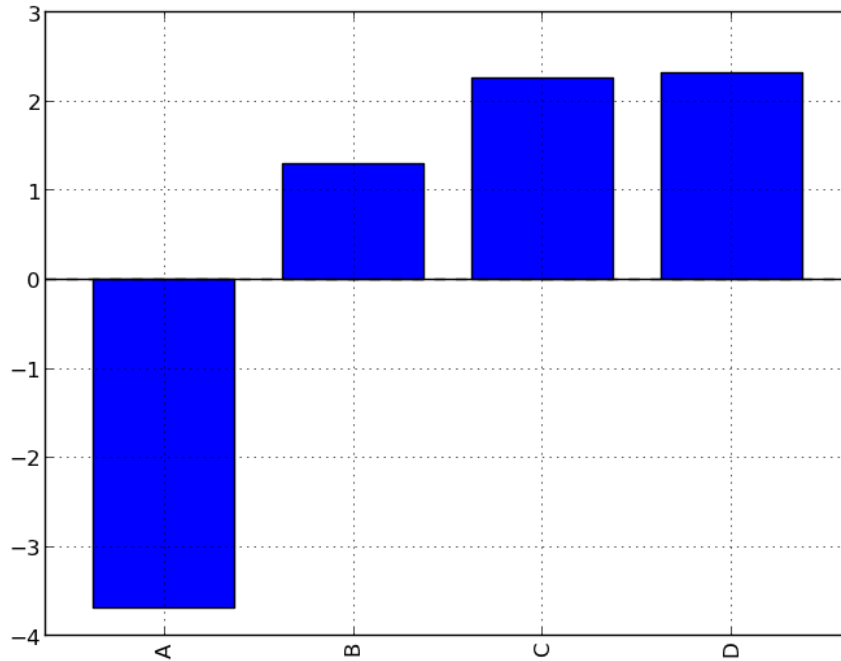
### 14.2.1 Bar plots

For labeled, non-time series data, you may wish to produce a bar plot:

```
In [1244]: plt.figure();
```

```
In [1244]: df.ix[5].plot(kind='bar'); plt.axhline(0, color='k')
```

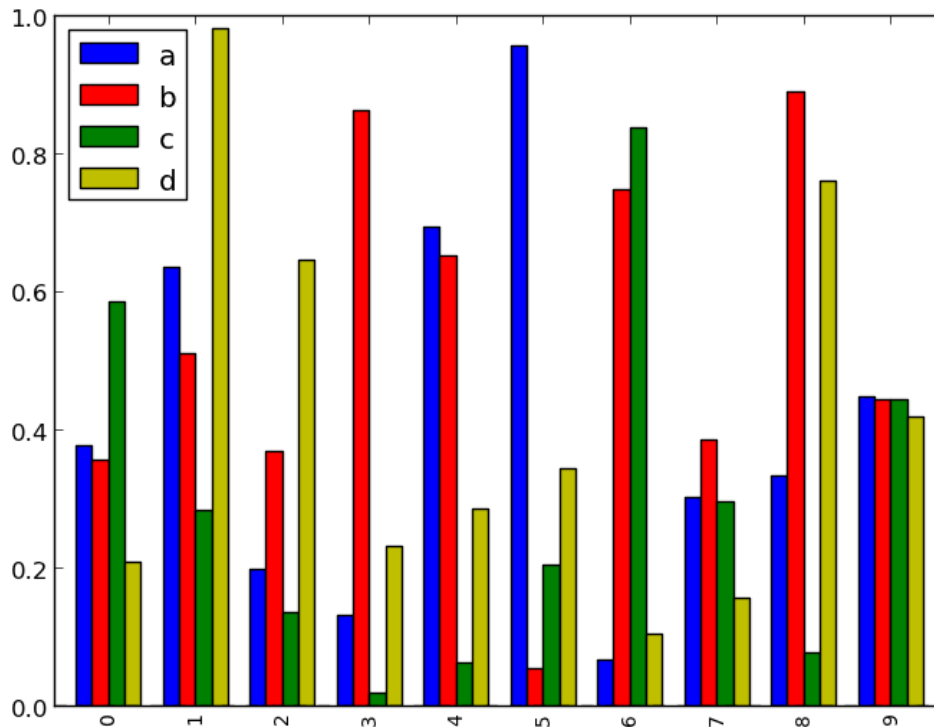
```
Out[1244]: <matplotlib.lines.Line2D at 0x127a89d0>
```



Calling a DataFrame's `plot` method with `kind='bar'` produces a multiple bar plot:

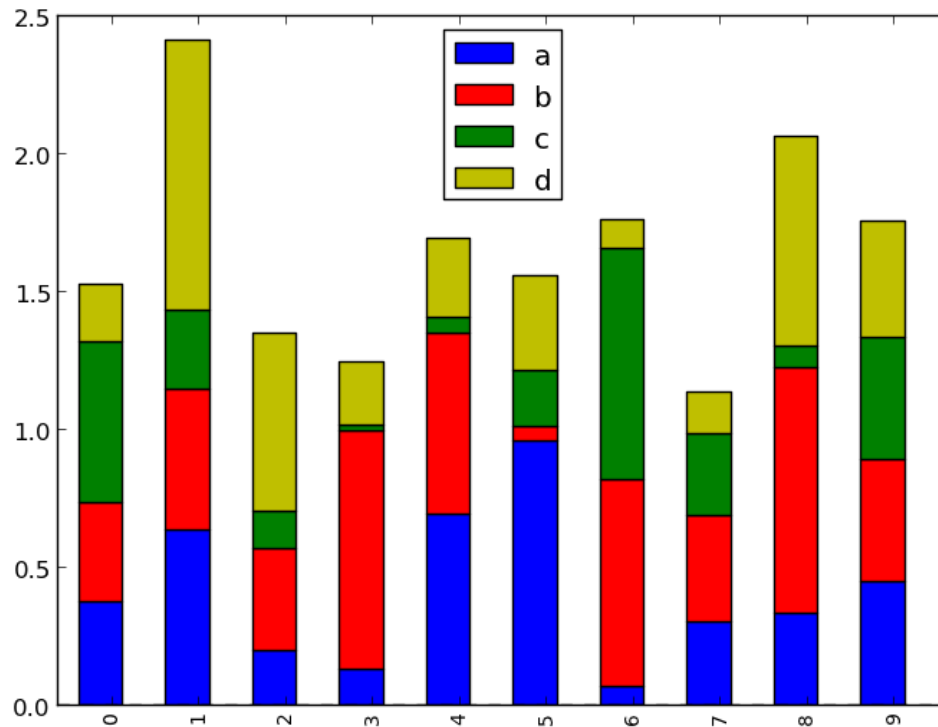
```
In [1245]: df2 = DataFrame(np.random.rand(10, 4), columns=['a', 'b', 'c', 'd'])
```

```
In [1246]: df2.plot(kind='bar');
```



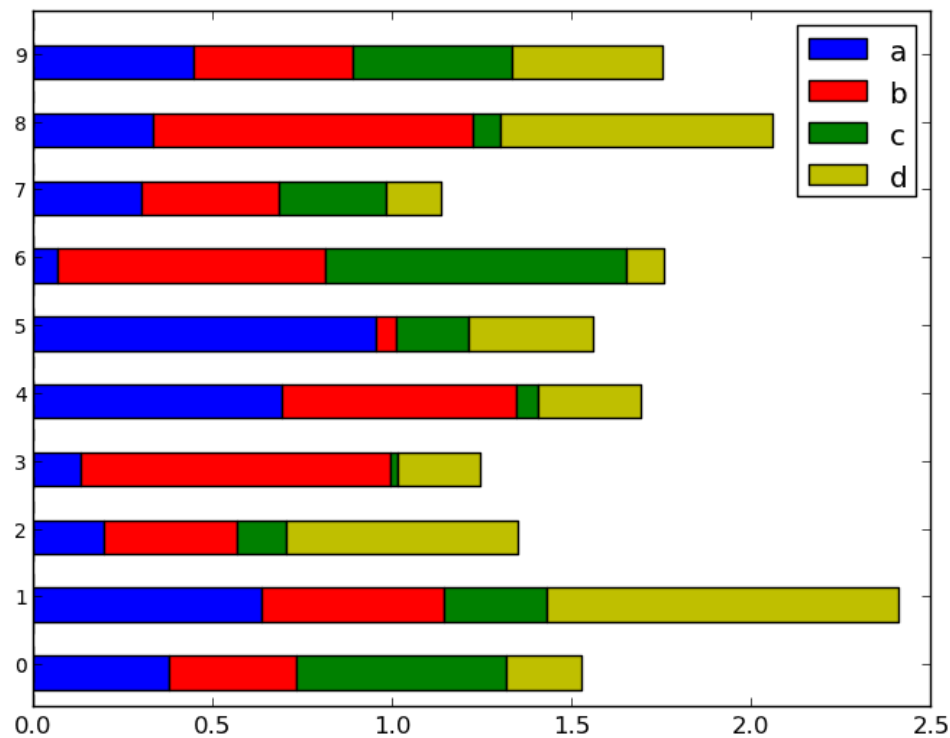
To produce a stacked bar plot, pass `stacked=True`:

```
In [1246]: df2.plot(kind='bar', stacked=True);
```



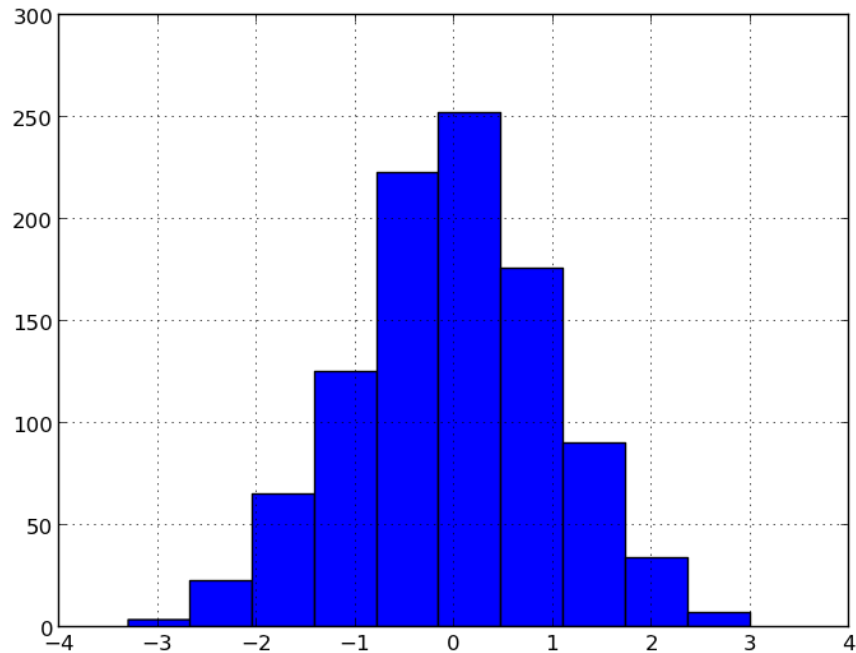
To get horizontal bar plots, pass `kind='barh'`:

```
In [1246]: df2.plot(kind='barh', stacked=True);
```



## 14.2.2 Histograms

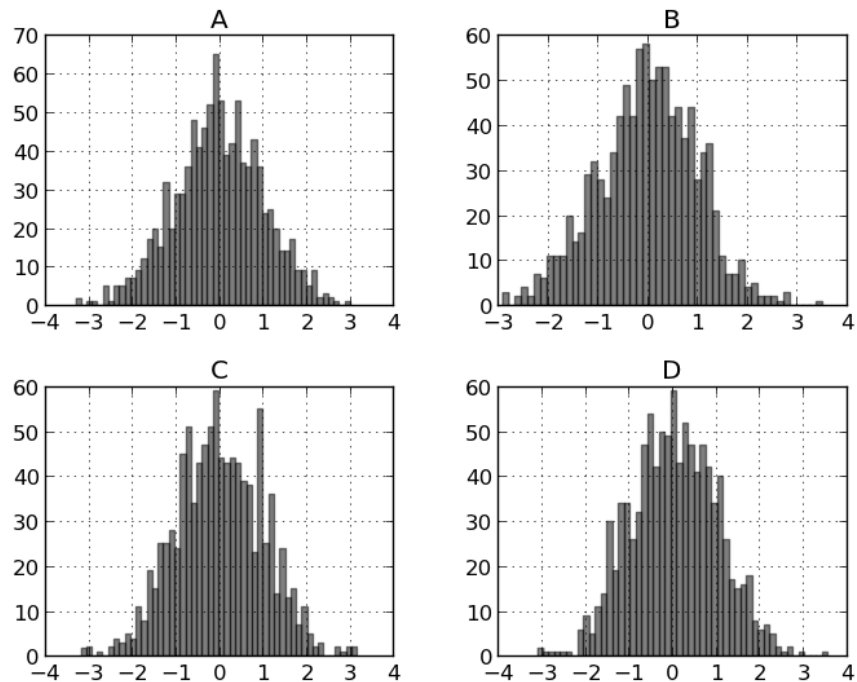
```
In [1246]: plt.figure();
In [1246]: df['A'].diff().hist()
Out[1246]: <matplotlib.axes.AxesSubplot at 0x136e8dd0>
```



For a DataFrame, hist plots the histograms of the columns on multiple subplots:

```
In [1247]: plt.figure()
Out[1247]: <matplotlib.figure.Figure at 0x1307aa10>

In [1248]: df.diff().hist(color='k', alpha=0.5, bins=50)
Out[1248]:
array([[Axes(0.125,0.552174;0.336957x0.347826),
        Axes(0.563043,0.552174;0.336957x0.347826)],
       [Axes(0.125,0.1;0.336957x0.347826),
        Axes(0.563043,0.1;0.336957x0.347826)]], dtype=object)
```



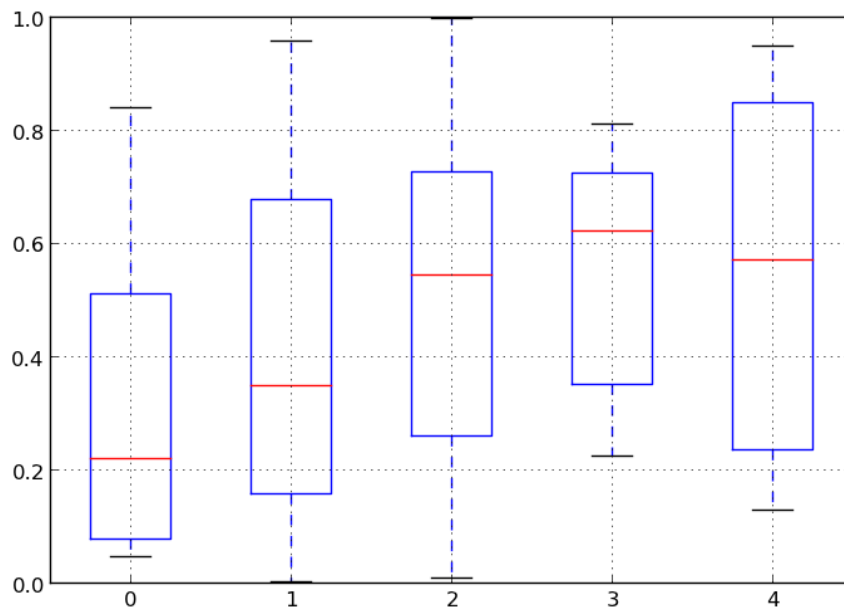
### 14.2.3 Box-Plotting

DataFrame has a `boxplot` method which allows you to visualize the distribution of values within each column. For instance, here is a boxplot representing five trials of 10 observations of a uniform random variable on  $[0,1]$ .

```
In [1249]: df = DataFrame(np.random.rand(10,5))
```

```
In [1250]: plt.figure();
```

```
In [1250]: bp = df.boxplot()
```



You can create a stratified boxplot using the `by` keyword argument to create groupings. For instance,

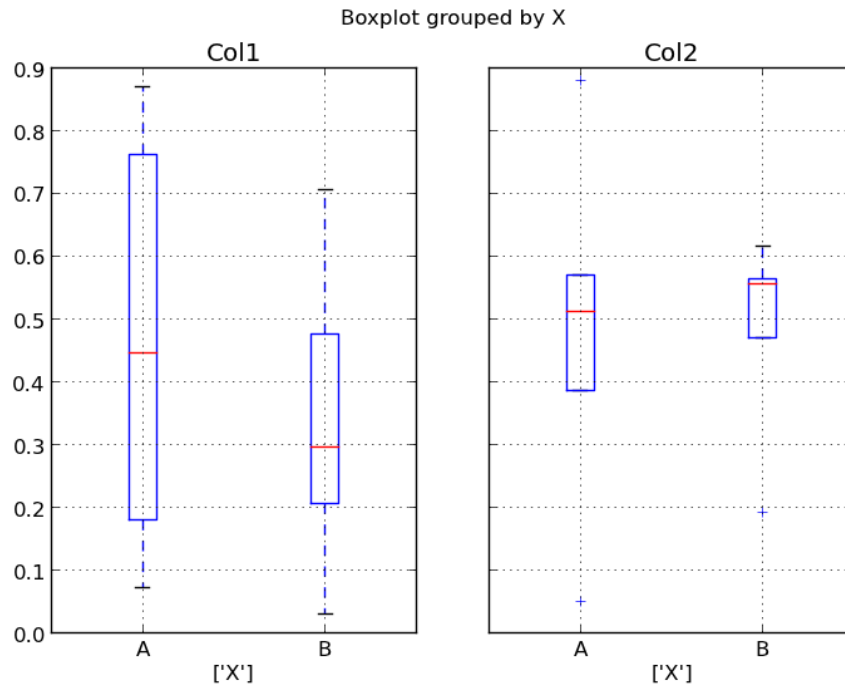


```
In [1251]: df = DataFrame(np.random.rand(10,2), columns=['Col1', 'Col2'] )
```

```
In [1252]: df['X'] = Series(['A','A','A','A','A','B','B','B','B','B'])
```

```
In [1253]: plt.figure();
```

```
In [1253]: bp = df.boxplot(by='X')
```



You can also pass a subset of columns to plot, as well as group by multiple columns:

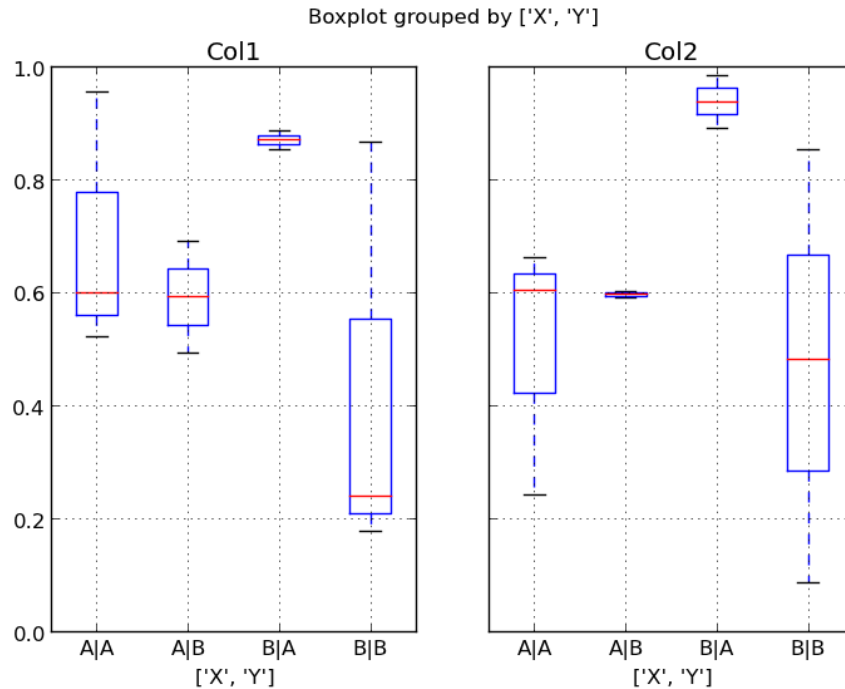
```
In [1254]: df = DataFrame(np.random.rand(10,3), columns=['Col1', 'Col2', 'Col3'])
```

```
In [1255]: df['X'] = Series(['A','A','A','A','A','B','B','B','B','B'])
```

```
In [1256]: df['Y'] = Series(['A','B','A','B','A','B','A','B','A','B'])
```

```
In [1257]: plt.figure();
```

```
In [1257]: bp = df.boxplot(column=['Col1','Col2'], by=['X','Y'])
```



## 14.2.4 Scatter plot matrix

**New in 0.7.3.** You can create a scatter plot matrix using the `scatter_matrix` method in `pandas.tools.plotting`:

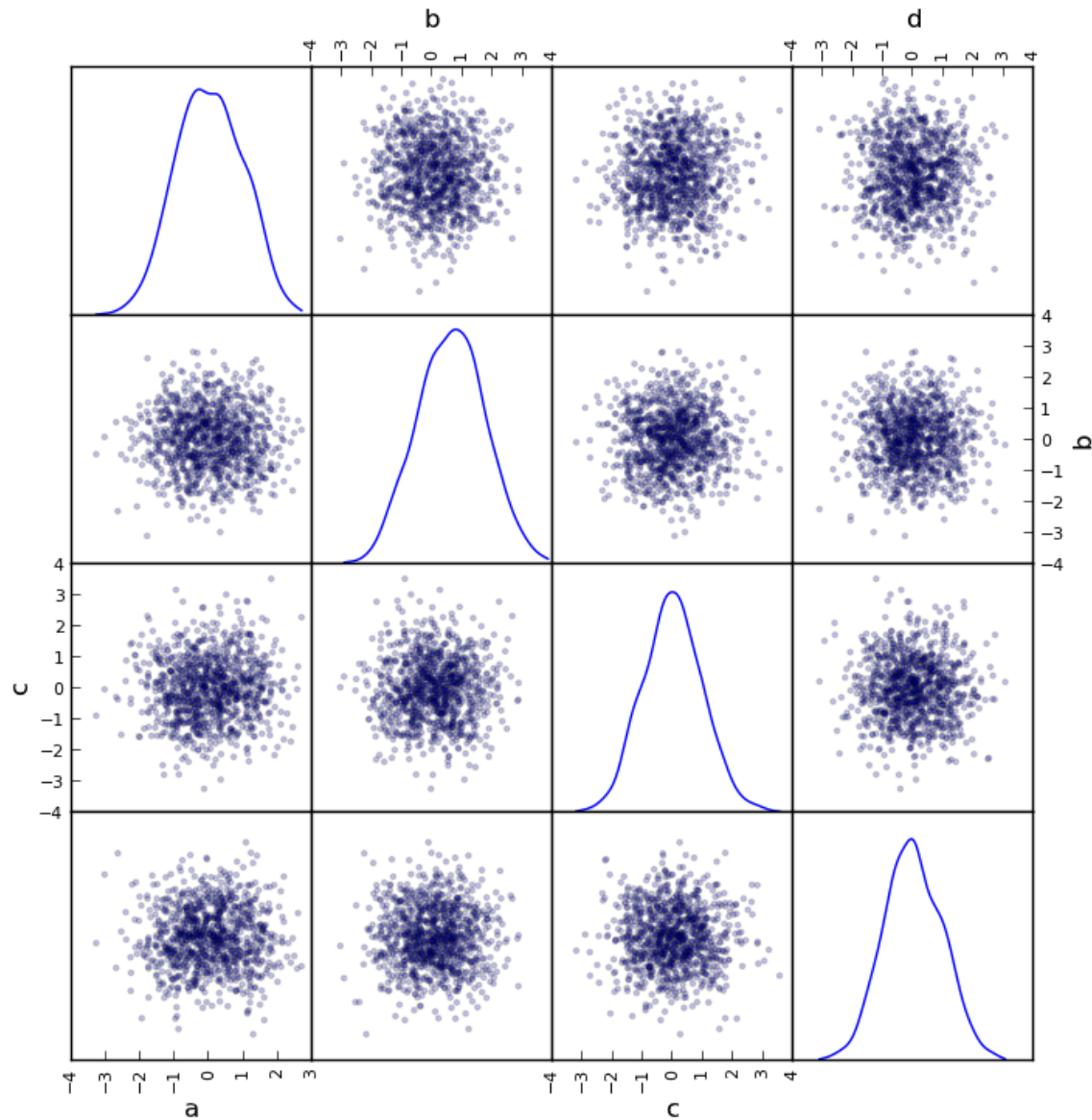
```
In [1258]: from pandas.tools.plotting import scatter_matrix
```

```
In [1259]: df = DataFrame(np.random.randn(1000, 4), columns=['a', 'b', 'c', 'd'])
```

```
In [1260]: scatter_matrix(df, alpha=0.2, figsize=(8, 8), diagonal='kde')
```

```
Out[1260]:
```

```
array([[Axes(0.125,0.7;0.19375x0.2), Axes(0.31875,0.7;0.19375x0.2),
       Axes(0.5125,0.7;0.19375x0.2), Axes(0.70625,0.7;0.19375x0.2)],
       [Axes(0.125,0.5;0.19375x0.2), Axes(0.31875,0.5;0.19375x0.2),
       Axes(0.5125,0.5;0.19375x0.2), Axes(0.70625,0.5;0.19375x0.2)],
       [Axes(0.125,0.3;0.19375x0.2), Axes(0.31875,0.3;0.19375x0.2),
       Axes(0.5125,0.3;0.19375x0.2), Axes(0.70625,0.3;0.19375x0.2)],
       [Axes(0.125,0.1;0.19375x0.2), Axes(0.31875,0.1;0.19375x0.2),
       Axes(0.5125,0.1;0.19375x0.2), Axes(0.70625,0.1;0.19375x0.2)]] , dtype=object)
```

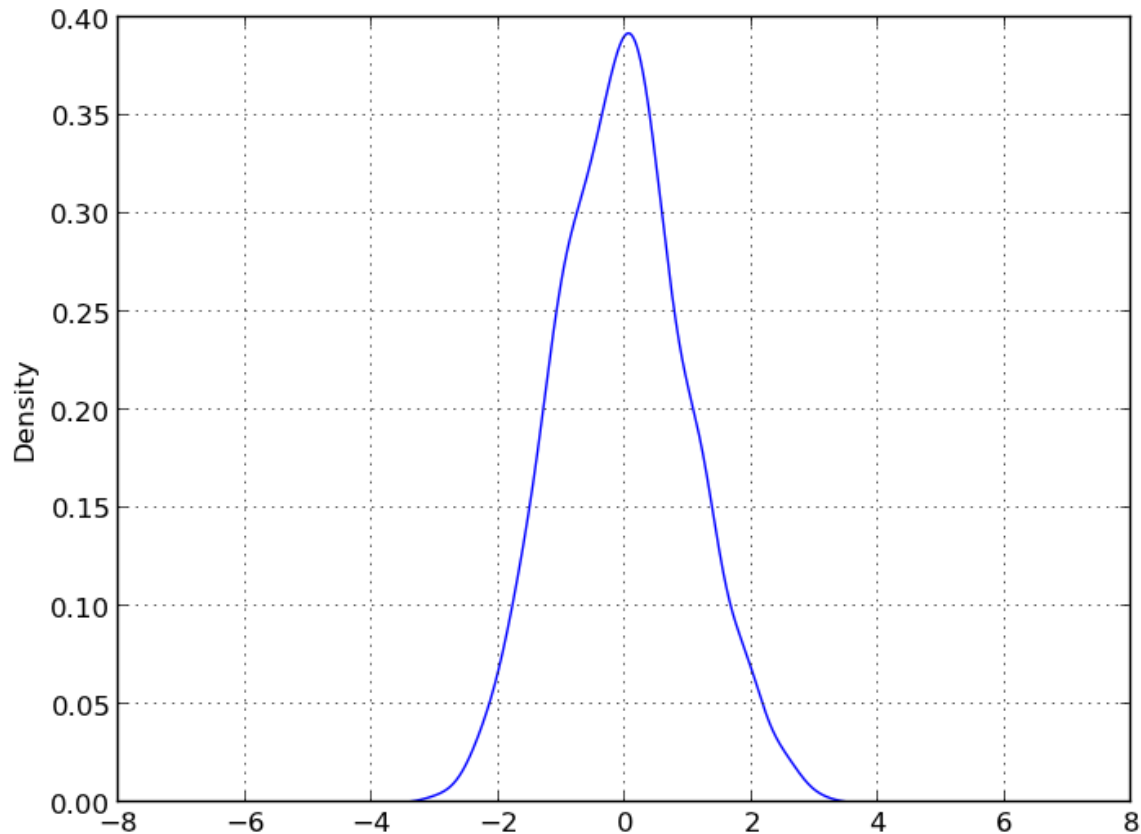


*New in*

0.8.0 You can create density plots using the `Series/DataFrame.plot` and setting `kind='kde'`:

```
In [1261]: ser = Series(np.random.randn(1000))

In [1262]: ser.plot(kind='kde')
Out[1262]: <matplotlib.axes.AxesSubplot at 0x16a11110>
```



### 14.2.5 Andrews Curves

Andrews curves allow one to plot multivariate data as a large number of curves that are created using the attributes of samples as coefficients for Fourier series. By coloring these curves differently for each class it is possible to visualize data clustering. Curves belonging to samples of the same class will usually be closer together and form larger structures.

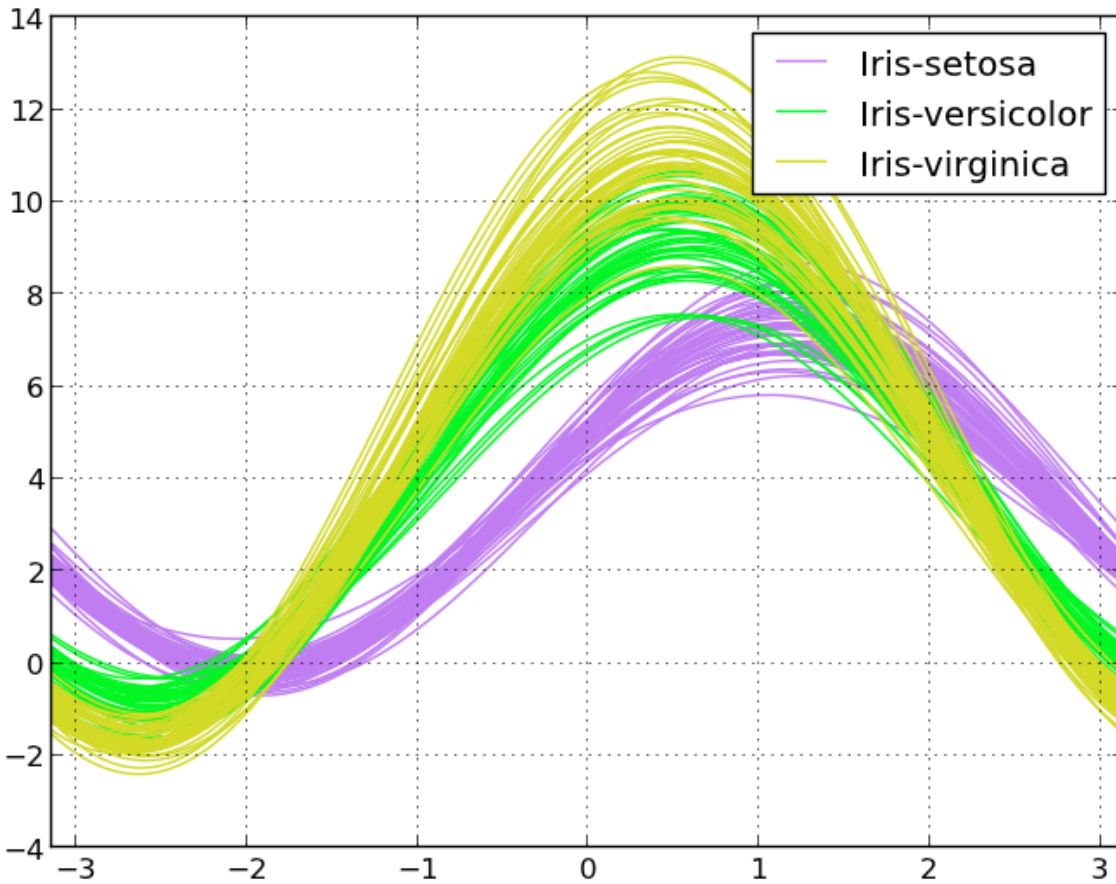
```
In [1263]: from pandas import read_csv

In [1264]: from pandas.tools.plotting import andrews_curves

In [1265]: data = read_csv('data/iris.data')

In [1266]: plt.figure()
Out[1266]: <matplotlib.figure.Figure at 0x16b58ad0>

In [1267]: andrews_curves(data, 'Name')
Out[1267]: <matplotlib.axes.AxesSubplot at 0x16b64c10>
```



### 14.2.6 Parallel Coordinates

Parallel coordinates is a plotting technique for plotting multivariate data. It allows one to see clusters in data and to estimate other statistics visually. Using parallel coordinates points are represented as connected line segments. Each vertical line represents one attribute. One set of connected line segments represents one data point. Points that tend to cluster will appear closer together.

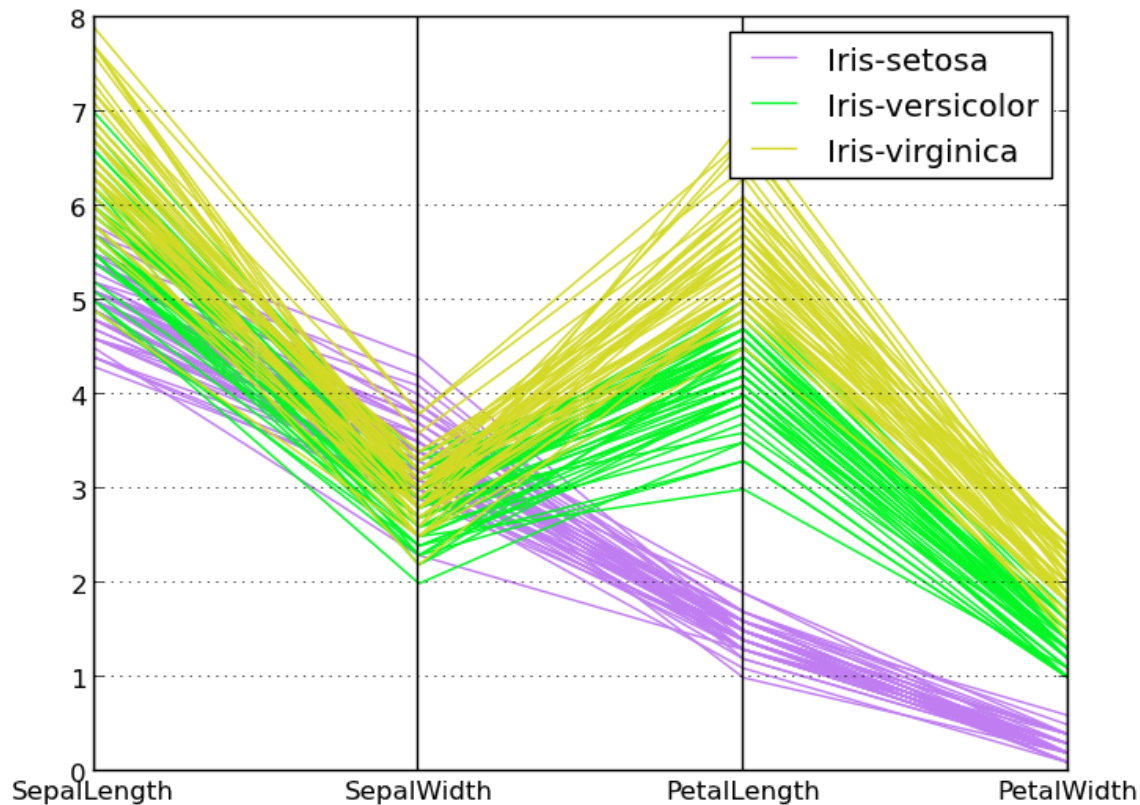
```
In [1268]: from pandas import read_csv

In [1269]: from pandas.tools.plotting import parallel_coordinates

In [1270]: data = read_csv('data/iris.data')

In [1271]: plt.figure()
Out[1271]: <matplotlib.figure.Figure at 0x16b589d0>

In [1272]: parallel_coordinates(data, 'Name')
Out[1272]: <matplotlib.axes.AxesSubplot at 0x1769f450>
```



### 14.2.7 Lag Plot

Lag plots are used to check if a data set or time series is random. Random data should not exhibit any structure in the lag plot. Non-random structure implies that the underlying data are not random.

```
In [1273]: from pandas.tools.plotting import lag_plot
```

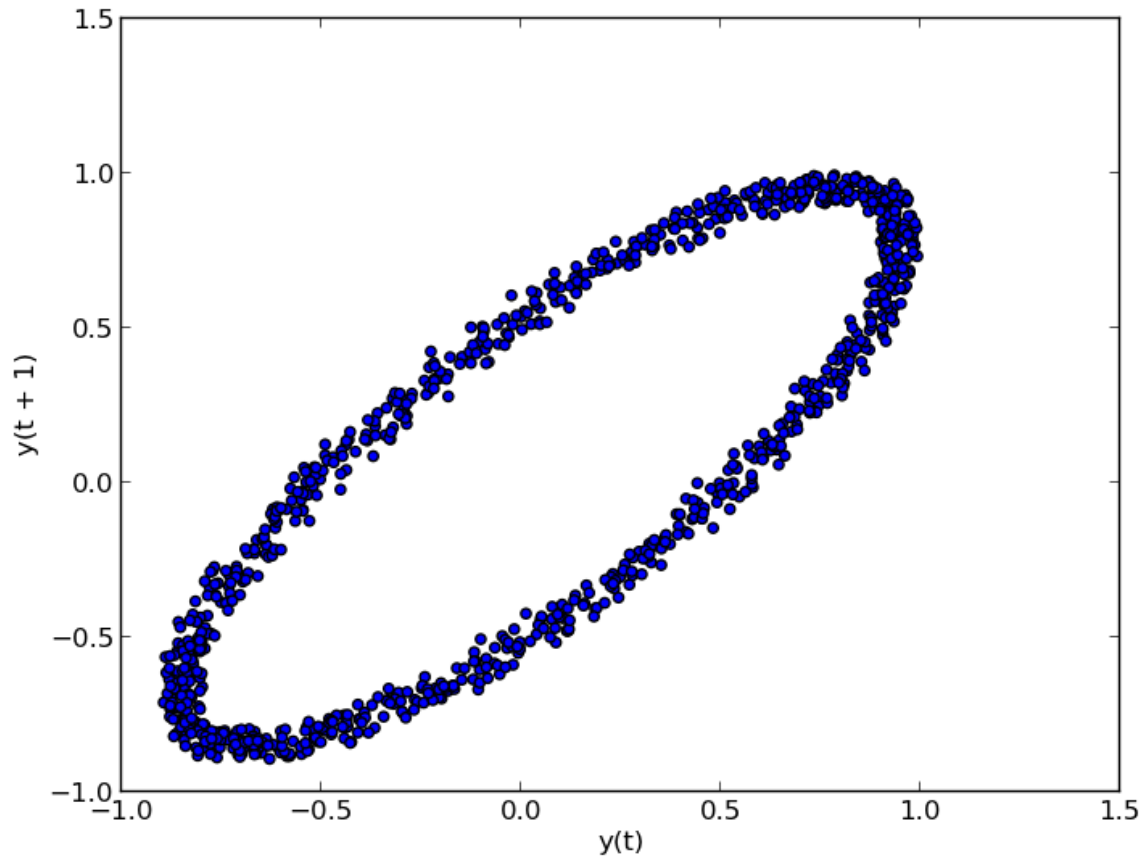
```
In [1274]: plt.figure()
```

```
Out[1274]: <matplotlib.figure.Figure at 0x17e75210>
```

```
In [1275]: data = Series(0.1 * np.random.random(1000) +
.....:    0.9 * np.sin(np.linspace(-99 * np.pi, 99 * np.pi, num=1000)))
.....:
```

```
In [1276]: lag_plot(data)
```

```
Out[1276]: <matplotlib.axes.AxesSubplot at 0x17e5a1d0>
```



### 14.2.8 Autocorrelation Plot

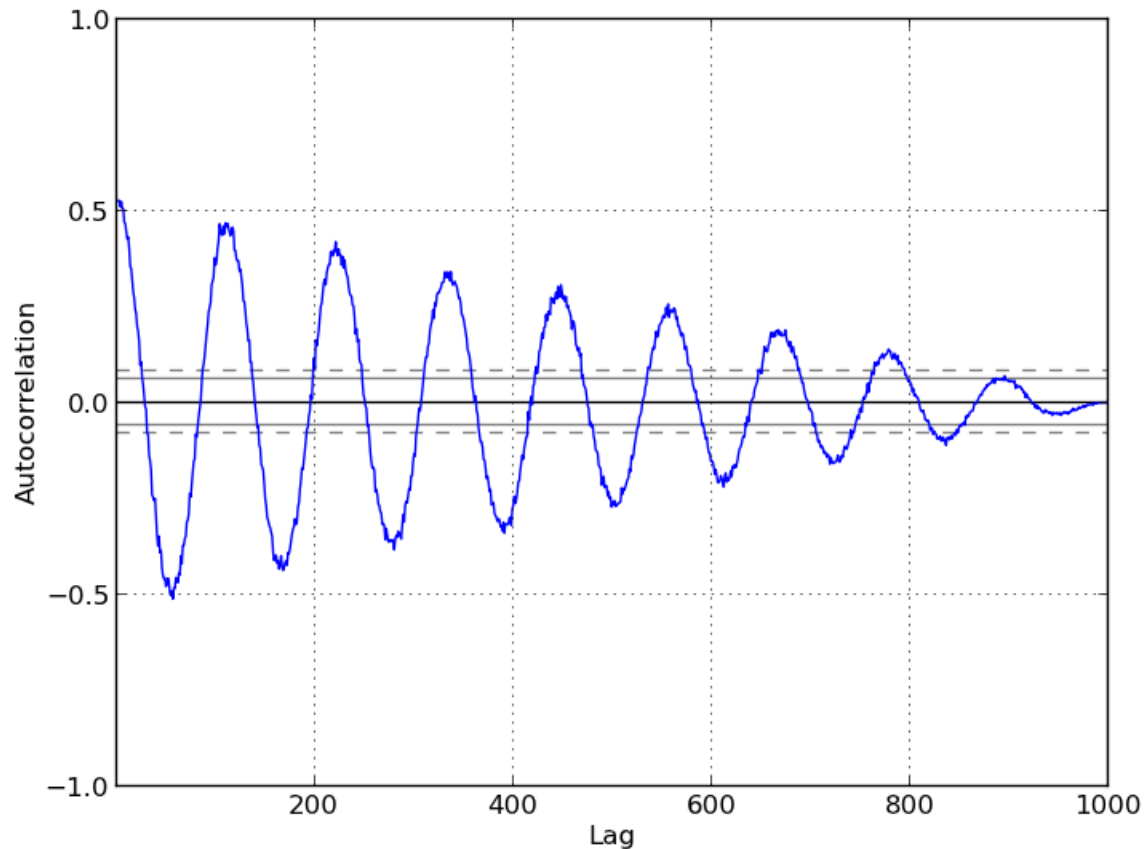
Autocorrelation plots are often used for checking randomness in time series. This is done by computing autocorrelations for data values at varying time lags. If time series is random, such autocorrelations should be near zero for any and all time-lag separations. If time series is non-random then one or more of the autocorrelations will be significantly non-zero. The horizontal lines displayed in the plot correspond to 95% and 99% confidence bands. The dashed line is 99% confidence band.

```
In [1277]: from pandas.tools.plotting import autocorrelation_plot

In [1278]: plt.figure()
Out[1278]: <matplotlib.figure.Figure at 0x17febe10>

In [1279]: data = Series(0.7 * np.random.random(1000) +
.....:     0.3 * np.sin(np.linspace(-9 * np.pi, 9 * np.pi, num=1000)))
.....:

In [1280]: autocorrelation_plot(data)
Out[1280]: <matplotlib.axes.AxesSubplot at 0x18004c10>
```



### 14.2.9 Bootstrap Plot

Bootstrap plots are used to visually assess the uncertainty of a statistic, such as mean, median, midrange, etc. A random subset of a specified size is selected from a data set, the statistic in question is computed for this subset and the process is repeated a specified number of times. Resulting plots and histograms are what constitutes the bootstrap plot.

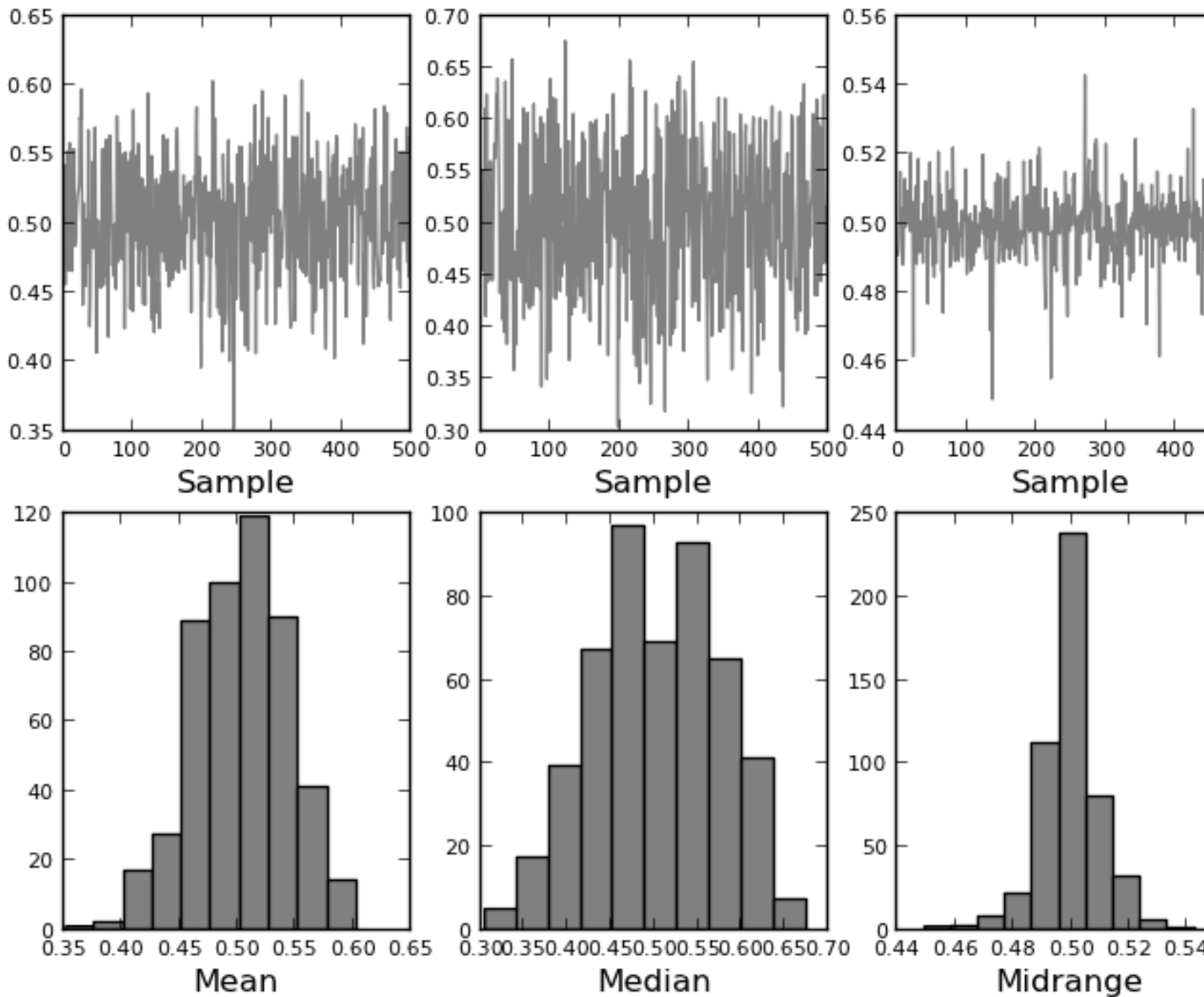
```
In [1281]: from pandas.tools.plotting import bootstrap_plot
```

```
In [1282]: data = Series(np.random.random(1000))
```

```
In [1283]: bootstrap_plot(data, size=50, samples=500, color='grey')
```

```
Out[1283]: <matplotlib.figure.Figure at 0x1769a3d0>
```





### 14.2.10 RadViz

RadViz is a way of visualizing multi-variate data. It is based on a simple spring tension minimization algorithm. Basically you set up a bunch of points in a plane. In our case they are equally spaced on a unit circle. Each point represents a single attribute. You then pretend that each sample in the data set is attached to each of these points by a spring, the stiffness of which is proportional to the numerical value of that attribute (they are normalized to unit interval). The point in the plane, where our sample settles to (where the forces acting on our sample are at an equilibrium) is where a dot representing our sample will be drawn. Depending on which class that sample belongs to it will be colored differently.

```
In [1284]: from pandas import read_csv
```

```
In [1285]: from pandas.tools.plotting import radviz
```

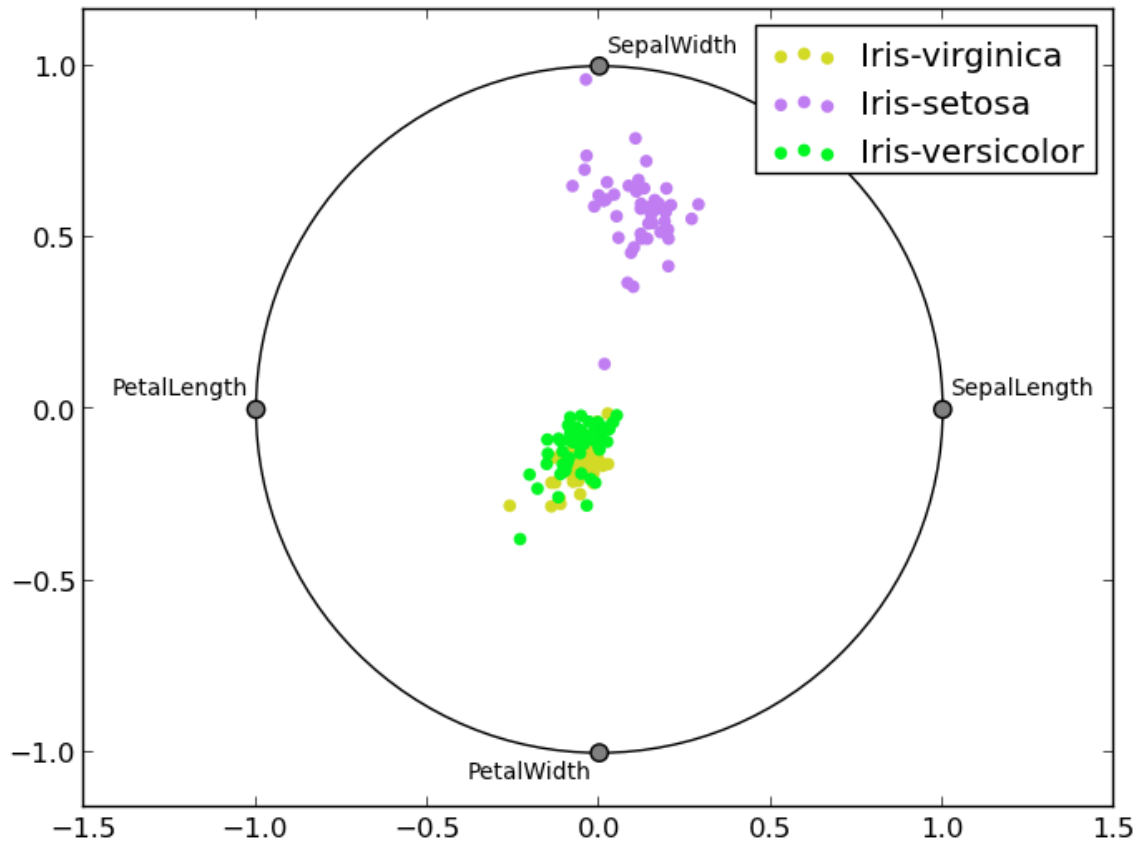
```
In [1286]: data = read_csv('data/iris.data')
```

```
In [1287]: plt.figure()
```

```
Out[1287]: <matplotlib.figure.Figure at 0x18437590>
```

```
In [1288]: radviz(data, 'Name')
```

```
Out[1288]: <matplotlib.axes.AxesSubplot at 0x194e5d90>
```



# IO TOOLS (TEXT, CSV, HDF5, ...)

## 15.1 Clipboard

A handy way to grab data is to use the `read_clipboard` method, which takes the contents of the clipboard buffer and passes them to the `read_table` method described in the next section. For instance, you can copy the following text to the clipboard (CTRL-C on many operating systems):

```
A B C
x 1 4 p
y 2 5 q
z 3 6 r
```

And then import the data directly to a `DataFrame` by calling:

```
clipdf = read_clipboard(sep='\s*')
```

```
In [738]: clipdf
```

```
Out[738]:
```

```
   A  B  C
x  1  4  p
y  2  5  q
z  3  6  r
```

## 15.2 CSV & Text files

The two workhorse functions for reading text files (a.k.a. flat files) are `read_csv()` and `read_table()`. They both use the same parsing code to intelligently convert tabular data into a `DataFrame` object. They can take a number of arguments:

- `filepath_or_buffer`: Either a string path to a file, or any object with a `read` method (such as an open file or `StringIO`).
- `sep` or `delimiter`: A delimiter / separator to split fields on. `read_csv` is capable of inferring the delimiter automatically in some cases by “sniffing.” The separator may be specified as a regular expression; for instance you may use `'s*'` to indicate arbitrary whitespace.
- `dialect`: string or `csv.Dialect` instance to expose more ways to specify the file format
- `header`: row number to use as the column names, and the start of the data. Defaults to 0 (first row); specify `None` if there is no header row.
- `skiprows`: A collection of numbers for rows in the file to skip. Can also be an integer to skip the first `n` rows

- `index_col`: column number, column name, or list of column numbers/names, to use as the `index` (row labels) of the resulting `DataFrame`. By default, it will number the rows without using any column, unless there is one more data column than there are headers, in which case the first column is taken as the index.
- `names`: List of column names to use. If passed, header will be implicitly set to `None`.
- `na_values`: optional list of strings to recognize as `NaN` (missing values), in addition to a default set. If you pass an empty list or an empty list for a particular column, no values (including empty strings) will be considered `NA`
- `parse_dates`: if `True` then index will be parsed as dates (`False` by default). You can specify more complicated options to parse a subset of columns or a combination of columns into a single date column (list of ints or names, list of lists, or dict) `[1, 2, 3]` -> try parsing columns 1, 2, 3 each as a separate date column `[[1, 3]]` -> combine columns 1 and 3 and parse as a single date column `{ 'foo' : [1, 3] }` -> parse columns 1, 3 as date and call result 'foo'
- `keep_date_col`: if `True`, then date component columns passed into `parse_dates` will be retained in the output (`False` by default).
- `date_parser`: function to use to parse strings into datetime objects. If `parse_dates` is `True`, it defaults to the very robust `dateutil.parser`. Specifying this implicitly sets `parse_dates` as `True`. You can also use functions from community supported date converters from `date_converters.py`
- `dayfirst`: if `True` then uses the `DD/MM` international/European date format (This is `False` by default)
- `thousands`: specifies the thousands separator. If not `None`, then parser will try to look for it in the output and parse relevant data to integers. Because it has to essentially scan through the data again, this causes a significant performance hit so only use if necessary.
- `comment`: denotes the start of a comment and ignores the rest of the line. Currently line commenting is not supported.
- `nrows`: Number of rows to read out of the file. Useful to only read a small portion of a large file
- `iterator`: If `True`, return a `TextParser` to enable reading a file into memory piece by piece
- `chunksize`: An number of rows to be used to “chunk” a file into pieces. Will cause an `TextParser` object to be returned. More on this below in the section on *iterating and chunking*
- `skip_footer`: number of lines to skip at bottom of file (default 0)
- `converters`: a dictionary of functions for converting values in certain columns, where keys are either integers or column labels
- `encoding`: a string representing the encoding to use if the contents are non-ascii
- `verbose`: show number of `NA` values inserted in non-numeric columns
- `squeeze`: if `True` then output with only one column is turned into `Series`

Consider a typical CSV file containing, in this case, some time series data:

```
In [739]: print open('foo.csv').read()
date,A,B,C
20090101,a,1,2
20090102,b,3,4
20090103,c,4,5
```

The default for `read_csv` is to create a `DataFrame` with simple numbered rows:

```
In [740]: read_csv('foo.csv')
Out[740]:
```

	date	A	B	C
0	20090101	a	1	2

```
1 20090102 b 3 4
2 20090103 c 4 5
```

In the case of indexed data, you can pass the column number or column name you wish to use as the index:

```
In [741]: read_csv('foo.csv', index_col=0)
```

```
Out[741]:
```

	A	B	C
date			
20090101	a	1	2
20090102	b	3	4
20090103	c	4	5

```
In [742]: read_csv('foo.csv', index_col='date')
```

```
Out[742]:
```

	A	B	C
date			
20090101	a	1	2
20090102	b	3	4
20090103	c	4	5

You can also use a list of columns to create a hierarchical index:

```
In [743]: read_csv('foo.csv', index_col=[0, 'A'])
```

```
Out[743]:
```

		B	C
date	A		
20090101	a	1	2
20090102	b	3	4
20090103	c	4	5

The `dialect` keyword gives greater flexibility in specifying the file format. By default it uses the Excel dialect but you can specify either the dialect name or a `csv.Dialect` instance.

Suppose you had data with unenclosed quotes:

```
In [744]: print data
label1,label2,label3
index1,"a,c,e
index2,b,d,f
```

By default, `read_csv` uses the Excel dialect and treats the double quote as the quote character, which causes it to fail when it finds a newline before it finds the closing double quote.

We can get around this using `dialect`

```
In [745]: dia = csv.excel()
```

```
In [746]: dia.quoting = csv.QUOTE_NONE
```

```
In [747]: read_csv(StringIO(data), dialect=dia)
```

```
Out[747]:
```

	label1	label2	label3
index1	"a	c	e
index2	b	d	f

The parsers make every attempt to “do the right thing” and not be very fragile. Type inference is a pretty big deal. So if a column can be coerced to integer dtype without altering the contents, it will do so. Any non-numeric columns will come through as object dtype as with the rest of pandas objects.

### 15.2.1 Specifying Date Columns

To better facilitate working with datetime data, `read_csv()` and `read_table()` uses the keyword arguments `parse_dates` and `date_parser` to allow users to specify a variety of columns and date/time formats to turn the input text data into datetime objects.

The simplest case is to just pass in `parse_dates=True`:

```
# Use a column as an index, and parse it as dates.  
In [748]: df = read_csv('foo.csv', index_col=0, parse_dates=True)
```

```
In [749]: df  
Out[749]:
```

	A	B	C
date			
2009-01-01	a	1	2
2009-01-02	b	3	4
2009-01-03	c	4	5

```
# These are python datetime objects  
In [750]: df.index  
Out[750]:  
<class 'pandas.tseries.index.DatetimeIndex'>  
[2009-01-01 00:00:00, ..., 2009-01-03 00:00:00]  
Length: 3, Freq: None, Timezone: None
```

It is often the case that we may want to store date and time data separately, or store various date fields separately. the `parse_dates` keyword can be used to specify a combination of columns to parse the dates and/or times from.

You can specify a list of column lists to `parse_dates`, the resulting date columns will be prepended to the output (so as to not affect the existing column order) and the new column names will be the concatenation of the component column names:

```
In [751]: print open('tmp.csv').read()  
KORD,19990127, 19:00:00, 18:56:00, 0.8100  
KORD,19990127, 20:00:00, 19:56:00, 0.0100  
KORD,19990127, 21:00:00, 20:56:00, -0.5900  
KORD,19990127, 21:00:00, 21:18:00, -0.9900  
KORD,19990127, 22:00:00, 21:56:00, -0.5900  
KORD,19990127, 23:00:00, 22:56:00, -0.5900
```

```
In [752]: df = read_csv('tmp.csv', header=None, parse_dates=[[1, 2], [1, 3]])
```

```
In [753]: df  
Out[753]:
```

		X.2_X.3		X.2_X.4	X.1	X.5
0	1999-01-27	19:00:00	1999-01-27	18:56:00	KORD	0.81
1	1999-01-27	20:00:00	1999-01-27	19:56:00	KORD	0.01
2	1999-01-27	21:00:00	1999-01-27	20:56:00	KORD	-0.59
3	1999-01-27	21:00:00	1999-01-27	21:18:00	KORD	-0.99
4	1999-01-27	22:00:00	1999-01-27	21:56:00	KORD	-0.59
5	1999-01-27	23:00:00	1999-01-27	22:56:00	KORD	-0.59

By default the parser removes the component date columns, but you can choose to retain them via the `keep_date_col` keyword:

```
In [754]: df = read_csv('tmp.csv', header=None, parse_dates=[[1, 2], [1, 3]],  
.....:                  keep_date_col=True)  
.....:
```

```
In [755]: df
Out[755]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 6 entries, 0 to 5
Data columns:
X.2_X.3      6  non-null values
X.2_X.4      6  non-null values
X.1          6  non-null values
X.2          6  non-null values
X.3          6  non-null values
X.4          6  non-null values
X.5          6  non-null values
dtypes: float64(1), int64(1), object(5)
```

Note that if you wish to combine multiple columns into a single date column, a nested list must be used. In other words, `parse_dates=[1, 2]` indicates that the second and third columns should each be parsed as separate date columns while `parse_dates=[[1, 2]]` means the two columns should be parsed into a single column.

You can also use a dict to specify custom name columns:

```
In [756]: date_spec = {'nominal': [1, 2], 'actual': [1, 3]}

In [757]: df = read_csv('tmp.csv', header=None, parse_dates=date_spec)
```

```
In [758]: df
Out[758]:
```

	nominal		actual		X.1	X.5
0	1999-01-27	19:00:00	1999-01-27	18:56:00	KORD	0.81
1	1999-01-27	20:00:00	1999-01-27	19:56:00	KORD	0.01
2	1999-01-27	21:00:00	1999-01-27	20:56:00	KORD	-0.59
3	1999-01-27	21:00:00	1999-01-27	21:18:00	KORD	-0.99
4	1999-01-27	22:00:00	1999-01-27	21:56:00	KORD	-0.59
5	1999-01-27	23:00:00	1999-01-27	22:56:00	KORD	-0.59

It is important to remember that if multiple text columns are to be parsed into a single date column, then a new column is prepended to the data. The `index_col` specification is based off of this new set of columns rather than the original data columns:

```
In [759]: date_spec = {'nominal': [1, 2], 'actual': [1, 3]}

In [760]: df = read_csv('tmp.csv', header=None, parse_dates=date_spec,
.....:                  index_col=0) #index is the nominal column
.....:
```

```
In [761]: df
Out[761]:
```

	nominal		actual	X.1	X.5
1999-01-27	19:00:00	1999-01-27	18:56:00	KORD	0.81
1999-01-27	20:00:00	1999-01-27	19:56:00	KORD	0.01
1999-01-27	21:00:00	1999-01-27	20:56:00	KORD	-0.59
1999-01-27	21:00:00	1999-01-27	21:18:00	KORD	-0.99
1999-01-27	22:00:00	1999-01-27	21:56:00	KORD	-0.59
1999-01-27	23:00:00	1999-01-27	22:56:00	KORD	-0.59

## 15.2.2 Date Parsing Functions

Finally, the parser allows you can specify a custom `date_parser` function to take full advantage of the flexibility of the date parsing API:

```
In [762]: import pandas.io.date_converters as conv

In [763]: df = read_csv('tmp.csv', header=None, parse_dates=date_spec,
.....:                  date_parser=conv.parse_date_time)
.....:
```

```
In [764]: df
Out[764]:
```

	nominal	actual	X.1	X.5
0	1999-01-27 19:00:00	1999-01-27 18:56:00	KORD	0.81
1	1999-01-27 20:00:00	1999-01-27 19:56:00	KORD	0.01
2	1999-01-27 21:00:00	1999-01-27 20:56:00	KORD	-0.59
3	1999-01-27 21:00:00	1999-01-27 21:18:00	KORD	-0.99
4	1999-01-27 22:00:00	1999-01-27 21:56:00	KORD	-0.59
5	1999-01-27 23:00:00	1999-01-27 22:56:00	KORD	-0.59

You can explore the date parsing functionality in `date_converters.py` and add your own. We would love to turn this module into a community supported set of date/time parsers. To get you started, `date_converters.py` contains functions to parse dual date and time columns, year/month/day columns, and year/month/day/hour/minute/second columns. It also contains a `generic_parser` function so you can curry it with a function that deals with a single date rather than the entire array.

## 15.2.3 International Date Formats

While US date formats tend to be MM/DD/YYYY, many international formats use DD/MM/YYYY instead. For convenience, a `dayfirst` keyword is provided:

```
In [765]: print open('tmp.csv').read()
date,value,cat
1/6/2000,5,a
2/6/2000,10,b
3/6/2000,15,c

In [766]: read_csv('tmp.csv', parse_dates=[0])
Out[766]:
```

	date	value	cat
0	2000-01-06 00:00:00	5	a
1	2000-02-06 00:00:00	10	b
2	2000-03-06 00:00:00	15	c

```
In [767]: read_csv('tmp.csv', dayfirst=True, parse_dates=[0])
Out[767]:
```

	date	value	cat
0	2000-06-01 00:00:00	5	a
1	2000-06-02 00:00:00	10	b
2	2000-06-03 00:00:00	15	c

## 15.2.4 Thousand Separators

For large integers that have been written with a thousands separator, you can set the `thousands` keyword to `True` so that integers will be parsed correctly:



By default, integers with a thousands separator will be parsed as strings

```
In [768]: print open('tmp.csv').read()
ID|level|category
Patient1|123,000|x
Patient2|23,000|y
Patient3|1,234,018|z

In [769]: df = read_csv('tmp.csv', sep='|')
```

```
In [770]: df
Out[770]:
```

	ID	level	category
0	Patient1	123,000	x
1	Patient2	23,000	y
2	Patient3	1,234,018	z

```
In [771]: df.level.dtype
Out[771]: dtype('object')
```

The thousands keyword allows integers to be parsed correctly

```
In [772]: print open('tmp.csv').read()
ID|level|category
Patient1|123,000|x
Patient2|23,000|y
Patient3|1,234,018|z

In [773]: df = read_csv('tmp.csv', sep='|', thousands=',')
```

```
In [774]: df
Out[774]:
```

	ID	level	category
0	Patient1	123000	x
1	Patient2	23000	y
2	Patient3	1234018	z

```
In [775]: df.level.dtype
Out[775]: dtype('int64')
```

## 15.2.5 Comments

Sometimes comments or meta data may be included in a file:

```
In [776]: print open('tmp.csv').read()
ID,level,category
Patient1,123000,x # really unpleasant
Patient2,23000,y # wouldn't take his medicine
Patient3,1234018,z # awesome
```

By default, the parse includes the comments in the output:

```
In [777]: df = read_csv('tmp.csv')

In [778]: df
Out[778]:
```

	ID	level	category
0	Patient1	123000	x # really unpleasant

```
1 Patient2      23000  y # wouldn't take his medicine
2 Patient3     1234018      z # awesome
```

We can suppress the comments using the `comment` keyword:

```
In [779]: df = read_csv('tmp.csv', comment='#')
```

```
In [780]: df
```

```
Out[780]:
```

	ID	level	category
0	Patient1	123000	x
1	Patient2	23000	y
2	Patient3	1234018	z

## 15.2.6 Returning Series

Using the `squeeze` keyword, the parser will return output with a single column as a `Series`:

```
In [781]: print open('tmp.csv').read()
```

```
level
Patient1,123000
Patient2,23000
Patient3,1234018
```

```
In [782]: output = read_csv('tmp.csv', squeeze=True)
```

```
In [783]: output
```

```
Out[783]:
```

Patient1	123000
Patient2	23000
Patient3	1234018

Name: level

```
In [784]: type(output)
```

```
Out[784]: pandas.core.series.Series
```

## 15.2.7 Files with Fixed Width Columns

While `read_csv` reads delimited data, the `read_fwf()` function works with data files that have known and fixed column widths. The function parameters to `read_fwf` are largely the same as `read_csv` with two extra parameters:

- `colspecs`: a list of pairs (tuples), giving the extents of the fixed-width fields of each line as half-open intervals [from, to[
- `widths`: a list of field widths, which can be used instead of `colspecs` if the intervals are contiguous

Consider a typical fixed-width data file:

```
In [785]: print open('bar.csv').read()
```

```
id8141      360.242940      149.910199      11950.7
id1594      444.953632      166.985655      11788.4
id1849      364.136849      183.628767      11806.2
id1230      413.836124      184.375703      11916.8
id1948      502.953953      173.237159      12468.3
```

In order to parse this file into a `DataFrame`, we simply need to supply the column specifications to the `read_fwf` function along with the file name:

*#Column specifications are a list of half-intervals*

```
In [786]: colspecs = [(0, 6), (8, 20), (21, 33), (34, 43)]
```

```
In [787]: df = read_fwf('bar.csv', colspecs=colspecs, header=None, index_col=0)
```

```
In [788]: df
```

```
Out[788]:
```

	X.2	X.3	X.4
X.1			
id8141	360.242940	149.910199	11950.7
id1594	444.953632	166.985655	11788.4
id1849	364.136849	183.628767	11806.2
id1230	413.836124	184.375703	11916.8
id1948	502.953953	173.237159	12468.3

Note how the parser automatically picks column names X.<column number> when header=None argument is specified. Alternatively, you can supply just the column widths for contiguous columns:

*#Widths are a list of integers*

```
In [789]: widths = [6, 14, 13, 10]
```

```
In [790]: df = read_fwf('bar.csv', widths=widths, header=None)
```

```
In [791]: df
```

```
Out[791]:
```

	X.1	X.2	X.3	X.4
0	id8141	360.242940	149.910199	11950.7
1	id1594	444.953632	166.985655	11788.4
2	id1849	364.136849	183.628767	11806.2
3	id1230	413.836124	184.375703	11916.8
4	id1948	502.953953	173.237159	12468.3

The parser will take care of extra white spaces around the columns so it's ok to have extra separation between the columns in the file.

## 15.2.8 Files with an “implicit” index column

Consider a file with one less entry in the header than the number of data column:

```
In [792]: print open('foo.csv').read()
```

```
A,B,C
20090101,a,1,2
20090102,b,3,4
20090103,c,4,5
```

In this special case, read\_csv assumes that the first column is to be used as the index of the DataFrame:

```
In [793]: read_csv('foo.csv')
```

```
Out[793]:
```

	A	B	C
20090101	a	1	2
20090102	b	3	4
20090103	c	4	5

Note that the dates weren't automatically parsed. In that case you would need to do as before:

```
In [794]: df = read_csv('foo.csv', parse_dates=True)
```

```
In [795]: df.index
Out[795]:
<class 'pandas.tseries.index.DatetimeIndex'>
[2009-01-01 00:00:00, ..., 2009-01-03 00:00:00]
Length: 3, Freq: None, Timezone: None
```

## 15.2.9 Reading DataFrame objects with MultiIndex

Suppose you have data indexed by two columns:

```
In [796]: print open('data/minindex_ex.csv').read()
year,indiv,zit,xit
1977,"A",1.2,.6
1977,"B",1.5,.5
1977,"C",1.7,.8
1978,"A",.2,.06
1978,"B",.7,.2
1978,"C",.8,.3
1978,"D",.9,.5
1978,"E",1.4,.9
1979,"C",.2,.15
1979,"D",.14,.05
1979,"E",.5,.15
1979,"F",1.2,.5
1979,"G",3.4,1.9
1979,"H",5.4,2.7
1979,"I",6.4,1.2
```

The `index_col` argument to `read_csv` and `read_table` can take a list of column numbers to turn multiple columns into a MultiIndex:

```
In [797]: df = read_csv("data/minindex_ex.csv", index_col=[0,1])
```

```
In [798]: df
```

```
Out[798]:
```

		zit	xit
year	indiv		
1977	A	1.20	0.60
	B	1.50	0.50
	C	1.70	0.80
1978	A	0.20	0.06
	B	0.70	0.20
	C	0.80	0.30
	D	0.90	0.50
	E	1.40	0.90
1979	C	0.20	0.15
	D	0.14	0.05
	E	0.50	0.15
	F	1.20	0.50
	G	3.40	1.90
	H	5.40	2.70
	I	6.40	1.20

```
In [799]: df.ix[1978]
```

```
Out[799]:
```

	zit	xit
indiv		

```
A    0.2  0.06
B    0.7  0.20
C    0.8  0.30
D    0.9  0.50
E    1.4  0.90
```

### 15.2.10 Automatically “sniffing” the delimiter

`read_csv` is capable of inferring delimited (not necessarily comma-separated) files. YMMV, as pandas uses the `Sniffer` class of the `csv` module.

```
In [800]: print open('tmp2.csv').read()
:0:1:2:3
0:0.469112299907:-0.282863344329:-1.50905850317:-1.13563237102
1:1.21211202502:-0.173214649053:0.119208711297:-1.04423596628
2:-0.861848963348:-2.10456921889:-0.494929274069:1.07180380704
3:0.721555162244:-0.70677113363:-1.03957498511:0.271859885543
4:-0.424972329789:0.567020349794:0.276232019278:-1.08740069129
5:-0.673689708088:0.113648409689:-1.47842655244:0.524987667115
6:0.40470521868:0.57704598592:-1.71500201611:-1.03926848351
7:-0.370646858236:-1.15789225064:-1.34431181273:0.844885141425
8:1.07576978372:-0.10904997528:1.64356307036:-1.46938795954
9:0.357020564133:-0.67460010373:-1.77690371697:-0.968913812447
```

```
In [801]: read_csv('tmp2.csv')
Out[801]:
:0:1:2:3
0  0:0.469112299907:-0.282863344329:-1.50905850317:-
1  1:1.21211202502:-0.173214649053:0.119208711297:-1
2  2:-0.861848963348:-2.10456921889:-0.494929274069:
3  3:0.721555162244:-0.70677113363:-1.03957498511:0.
4  4:-0.424972329789:0.567020349794:0.276232019278:-
5  5:-0.673689708088:0.113648409689:-1.47842655244:0
6  6:0.40470521868:0.57704598592:-1.71500201611:-1.0
7  7:-0.370646858236:-1.15789225064:-1.34431181273:0
8  8:1.07576978372:-0.10904997528:1.64356307036:-1.4
9  9:0.357020564133:-0.67460010373:-1.77690371697:-0
```

### 15.2.11 Iterating through files chunk by chunk

Suppose you wish to iterate through a (potentially very large) file lazily rather than reading the entire file into memory, such as the following:

```
In [802]: print open('tmp.csv').read()
|0|1|2|3
0|0.469112299907|-0.282863344329|-1.50905850317|-1.13563237102
1|1.21211202502|-0.173214649053|0.119208711297|-1.04423596628
2|-0.861848963348|-2.10456921889|-0.494929274069|1.07180380704
3|0.721555162244|-0.70677113363|-1.03957498511|0.271859885543
4|-0.424972329789|0.567020349794|0.276232019278|-1.08740069129
5|-0.673689708088|0.113648409689|-1.47842655244|0.524987667115
6|0.40470521868|0.57704598592|-1.71500201611|-1.03926848351
7|-0.370646858236|-1.15789225064|-1.34431181273|0.844885141425
8|1.07576978372|-0.10904997528|1.64356307036|-1.46938795954
9|0.357020564133|-0.67460010373|-1.77690371697|-0.968913812447
```

```
In [803]: table = read_table('tmp.csv', sep='|')
```

```
In [804]: table
```

```
Out[804]:
```

	Unnamed: 0	0	1	2	3
0	0	0.469112	-0.282863	-1.509059	-1.135632
1	1	1.212112	-0.173215	0.119209	-1.044236
2	2	-0.861849	-2.104569	-0.494929	1.071804
3	3	0.721555	-0.706771	-1.039575	0.271860
4	4	-0.424972	0.567020	0.276232	-1.087401
5	5	-0.673690	0.113648	-1.478427	0.524988
6	6	0.404705	0.577046	-1.715002	-1.039268
7	7	-0.370647	-1.157892	-1.344312	0.844885
8	8	1.075770	-0.109050	1.643563	-1.469388
9	9	0.357021	-0.674600	-1.776904	-0.968914

By specifying a `chunksize` to `read_csv` or `read_table`, the return value will be an iterable object of type `TextParser`:

```
In [805]: reader = read_table('tmp.csv', sep='|', chunksize=4)
```

```
In [806]: reader
```

```
Out[806]: <pandas.io.parsers.TextParser at 0xa9ea350>
```

```
In [807]: for chunk in reader:
```

```
.....:     print chunk
```

```
.....:
```

```
Unamed: 0      0      1      2      3
0      0  0.469112 -0.282863 -1.509059 -1.135632
1      1  1.212112 -0.173215  0.119209 -1.044236
2      2 -0.861849 -2.104569 -0.494929  1.071804
3      3  0.721555 -0.706771 -1.039575  0.271860
Unamed: 0      0      1      2      3
0      4 -0.424972  0.567020  0.276232 -1.087401
1      5 -0.673690  0.113648 -1.478427  0.524988
2      6  0.404705  0.577046 -1.715002 -1.039268
3      7 -0.370647 -1.157892 -1.344312  0.844885
Unamed: 0      0      1      2      3
0      8  1.075770 -0.10905  1.643563 -1.469388
1      9  0.357021 -0.67460 -1.776904 -0.968914
```

Specifying `iterator=True` will also return the `TextParser` object:

```
In [808]: reader = read_table('tmp.csv', sep='|', iterator=True)
```

```
In [809]: reader.get_chunk(5)
```

```
Out[809]:
```

	Unnamed: 0	0	1	2	3
0	0	0.469112	-0.282863	-1.509059	-1.135632
1	1	1.212112	-0.173215	0.119209	-1.044236
2	2	-0.861849	-2.104569	-0.494929	1.071804
3	3	0.721555	-0.706771	-1.039575	0.271860
4	4	-0.424972	0.567020	0.276232	-1.087401

### 15.2.12 Writing to CSV format

The Series and DataFrame objects have an instance method `to_csv` which allows storing the contents of the object as a comma-separated-values file. The function takes a number of arguments. Only the first is required.

- `path`: A string path to the file to write
- `nanRep`: A string representation of a missing value (default `''`)
- `cols`: Columns to write (default `None`)
- `header`: Whether to write out the column names (default `True`)
- `index`: whether to write row (index) names (default `True`)
- `index_label`: Column label(s) for index column(s) if desired. If `None` (default), and `header` and `index` are `True`, then the index names are used. (A sequence should be given if the DataFrame uses `MultiIndex`).
- `mode`: Python write mode, default `'w'`
- `sep`: Field delimiter for the output file (default `','`)
- `encoding`: a string representing the encoding to use if the contents are non-ascii, for python versions prior to 3

### 15.2.13 Writing a formatted string

The DataFrame object has an instance method `to_string` which allows control over the string representation of the object. All arguments are optional:

- `buf` default `None`, for example a `StringIO` object
- `columns` default `None`, which columns to write
- `col_space` default `None`, number of spaces to write between columns
- `na_rep` default `NaN`, representation of NA value
- `formatters` default `None`, a dictionary (by column) of functions each of which takes a single argument and returns a formatted string
- `float_format` default `None`, a function which takes a single (float) argument and returns a formatted string; to be applied to floats in the DataFrame.
- `sparsify` default `True`, set to `False` for a DataFrame with a hierarchical index to print every multiindex key at each row.
- `index_names` default `True`, will print the names of the indices
- `index` default `True`, will print the index (ie, row labels)
- `header` default `True`, will print the column labels
- `justify` default `left`, will print column headers left- or right-justified

The Series object also has a `to_string` method, but with only the `buf`, `na_rep`, `float_format` arguments. There is also a `length` argument which, if set to `True`, will additionally output the length of the Series.

### 15.2.14 Writing to HTML format

DataFrame object has an instance method `to_html` which renders the contents of the DataFrame as an html table. The function arguments are as in the method `to_string` described above.

## 15.3 Excel files

The `ExcelFile` class can read an Excel 2003 file using the `xlrd` Python module and use the same parsing code as the above to convert tabular data into a `DataFrame`. To use it, create the `ExcelFile` object:

```
xls = ExcelFile('path_to_file.xls')
```

Then use the `parse` instance method with a `sheetname`, then use the same additional arguments as the parsers above:

```
xls.parse('Sheet1', index_col=None, na_values=['NA'])
```

To read sheets from an Excel 2007 file, you can pass a filename with a `.xlsx` extension, in which case the `openpyxl` module will be used to read the file.

It is often the case that users will insert columns to do temporary computations in Excel and you may not want to read in those columns. `ExcelFile.parse` takes a `parse_cols` keyword to allow you to specify a subset of columns to parse.

If `parse_cols` is an integer, then it is assumed to indicate the last column to be parsed.

```
xls.parse('Sheet1', parse_cols=2, index_col=None, na_values=['NA'])
```

If `parse_cols` is a list of integers, then it is assumed to be the file column indices to be parsed.

```
xls.parse('Sheet1', parse_cols=[0, 2, 3], index_col=None, na_values=['NA'])
```

To write a `DataFrame` object to a sheet of an Excel file, you can use the `to_excel` instance method. The arguments are largely the same as `to_csv` described above, the first argument being the name of the excel file, and the optional second argument the name of the sheet to which the `DataFrame` should be written. For example:

```
df.to_excel('path_to_file.xlsx', sheet_name='sheet1')
```

Files with a `.xls` extension will be written using `xlwt` and those with a `.xlsx` extension will be written using `openpyxl`. The `Panel` class also has a `to_excel` instance method, which writes each `DataFrame` in the `Panel` to a separate sheet.

In order to write separate `DataFrames` to separate sheets in a single Excel file, one can use the `ExcelWriter` class, as in the following example:

```
writer = ExcelWriter('path_to_file.xlsx')
df1.to_excel(writer, sheet_name='sheet1')
df2.to_excel(writer, sheet_name='sheet2')
writer.save()
```

## 15.4 HDF5 (PyTables)

`HDFStore` is a dict-like object which reads and writes pandas to the high performance HDF5 format using the excellent `PyTables` library.

```
In [810]: store = HDFStore('store.h5')
```

```
In [811]: print store
<class 'pandas.io.pytables.HDFStore'>
File path: store.h5
Empty
```

Objects can be written to the file just like adding key-value pairs to a dict:



```

In [812]: index = date_range('1/1/2000', periods=8)

In [813]: s = Series(randn(5), index=['a', 'b', 'c', 'd', 'e'])

In [814]: df = DataFrame(randn(8, 3), index=index,
.....:                  columns=['A', 'B', 'C'])
.....:

In [815]: wp = Panel(randn(2, 5, 4), items=['Item1', 'Item2'],
.....:                major_axis=date_range('1/1/2000', periods=5),
.....:                minor_axis=['A', 'B', 'C', 'D'])
.....:

In [816]: store['s'] = s

In [817]: store['df'] = df

In [818]: store['wp'] = wp

In [819]: store
Out[819]:
<class 'pandas.io.pytables.HDFStore'>
File path: store.h5
df      DataFrame
s       Series
wp      Panel

```

In a current or later Python session, you can retrieve stored objects:

```

In [820]: store['df']
Out[820]:
           A           B           C
2000-01-01 -0.362543 -0.006154 -0.923061
2000-01-02  0.895717  0.805244 -1.206412
2000-01-03  2.565646  1.431256  1.340309
2000-01-04 -1.170299 -0.226169  0.410835
2000-01-05  0.813850  0.132003 -0.827317
2000-01-06 -0.076467 -1.187678  1.130127
2000-01-07 -1.436737 -1.413681  1.607920
2000-01-08  1.024180  0.569605  0.875906

```



# SPARSE DATA STRUCTURES

We have implemented “sparse” versions of Series, DataFrame, and Panel. These are not sparse in the typical “mostly 0”. You can view these objects as being “compressed” where any data matching a specific value (NaN/missing by default, though any value can be chosen) is omitted. A special `SparseIndex` object tracks where data has been “sparsified”. This will make much more sense in an example. All of the standard pandas data structures have a `to_sparse` method:

```
In [1044]: ts = Series(randn(10))
```

```
In [1045]: ts[2:-2] = np.nan
```

```
In [1046]: sts = ts.to_sparse()
```

```
In [1047]: sts
```

```
Out[1047]:
```

```
0    0.469112
1   -0.282863
2         NaN
3         NaN
4         NaN
5         NaN
6         NaN
7         NaN
8   -0.861849
9   -2.104569
```

```
BlockIndex
```

```
Block locations: array([0, 8], dtype=int32)
```

```
Block lengths: array([2, 2], dtype=int32)
```

The `to_sparse` method takes a `kind` argument (for the sparse index, see below) and a `fill_value`. So if we had a mostly zero Series, we could convert it to sparse with `fill_value=0`:

```
In [1048]: ts.fillna(0).to_sparse(fill_value=0)
```

```
Out[1048]:
```

```
0    0.469112
1   -0.282863
2    0.000000
3    0.000000
4    0.000000
5    0.000000
6    0.000000
7    0.000000
8   -0.861849
9   -2.104569
```

```
BlockIndex
```

```
Block locations: array([0, 8], dtype=int32)
Block lengths: array([2, 2], dtype=int32)
```

The sparse objects exist for memory efficiency reasons. Suppose you had a large, mostly NA DataFrame:

```
In [1049]: df = DataFrame(randn(10000, 4))

In [1050]: df.ix[:9998] = np.nan

In [1051]: sdf = df.to_sparse()

In [1052]: sdf
Out[1052]:
<class 'pandas.sparse.frame.SparseDataFrame'>
Int64Index: 10000 entries, 0 to 9999
Columns: 4 entries, 0 to 3
dtypes: float64(4)
```

```
In [1053]: sdf.density
Out[1053]: 0.0001
```

As you can see, the density (% of values that have not been “compressed”) is extremely low. This sparse object takes up much less memory on disk (pickled) and in the Python interpreter. Functionally, their behavior should be nearly identical to their dense counterparts.

Any sparse object can be converted back to the standard dense form by calling `to_dense`:

```
In [1054]: sts.to_dense()
Out[1054]:
0    0.469112
1   -0.282863
2         NaN
3         NaN
4         NaN
5         NaN
6         NaN
7         NaN
8   -0.861849
9   -2.104569
```

## 16.1 SparseArray

`SparseArray` is the base layer for all of the sparse indexed data structures. It is a 1-dimensional ndarray-like object storing only values distinct from the `fill_value`:

```
In [1055]: arr = np.random.randn(10)

In [1056]: arr[2:5] = np.nan; arr[7:8] = np.nan

In [1057]: sparr = SparseArray(arr)

In [1058]: sparr
Out[1058]:
SparseArray([-1.9557, -1.6589,      nan,      nan,      nan,  1.1589,  0.1453,
            nan,  0.606 ,  1.3342])
IntIndex
Indices: array([0, 1, 5, 6, 8, 9], dtype=int32)
```

Like the indexed objects (SparseSeries, SparseDataFrame, SparsePanel), a SparseArray can be converted back to a regular ndarray by calling `to_dense`:

```
In [1059]: sparr.to_dense()
Out[1059]:
array([-1.9557, -1.6589,      nan,      nan,      nan,  1.1589,  0.1453,
        nan,  0.606 ,  1.3342])
```

## 16.2 SparseList

SparseList is a list-like data structure for managing a dynamic collection of SparseArrays. To create one, simply call the SparseList constructor with a `fill_value` (defaulting to NaN):

```
In [1060]: spl = SparseList()

In [1061]: spl
Out[1061]:
<pandas.sparse.list.SparseList object at 0xba580d0>
```

The two important methods are `append` and `to_array`. `append` can accept scalar values or any 1-dimensional sequence:

```
In [1062]: spl.append(np.array([1., nan, nan, 2., 3.]))

In [1063]: spl.append(5)

In [1064]: spl.append(sparr)

In [1065]: spl
Out[1065]:
<pandas.sparse.list.SparseList object at 0xba580d0>
SparseArray([ 1., nan, nan,  2.,  3.])
IntIndex
Indices: array([0, 3, 4], dtype=int32)
SparseArray([ 5.])
IntIndex
Indices: array([0], dtype=int32)
SparseArray([-1.9557, -1.6589,      nan,      nan,      nan,  1.1589,  0.1453,
        nan,  0.606 ,  1.3342])
IntIndex
Indices: array([0, 1, 5, 6, 8, 9], dtype=int32)
```

As you can see, all of the contents are stored internally as a list of memory-efficient SparseArray objects. Once you've accumulated all of the data, you can call `to_array` to get a single SparseArray with all the data:

```
In [1066]: spl.to_array()
Out[1066]:
SparseArray([ 1.      ,      nan,      nan,  2.      ,  3.      ,  5.      , -1.9557,
        -1.6589,      nan,      nan,      nan,  1.1589,  0.1453,      nan,
         0.606 ,  1.3342])
IntIndex
Indices: array([ 0,  3,  4,  5,  6,  7, 11, 12, 14, 15], dtype=int32)
```

## 16.3 SparseIndex objects

Two kinds of `SparseIndex` are implemented, `block` and `integer`. We recommend using `block` as it's more memory efficient. The `integer` format keeps an array of all of the locations where the data are not equal to the fill value. The `block` format tracks only the locations and sizes of blocks of data.

# CAVEATS AND GOTCHAS

## 17.1 NaN, Integer NA values and NA type promotions

### 17.1.1 Choice of NA representation

For lack of NA (missing) support from the ground up in NumPy and Python in general, we were given the difficult choice between either

- A *masked array* solution: an array of data and an array of boolean values indicating whether a value
- Using a special sentinel value, bit pattern, or set of sentinel values to denote NA across the dtypes

For many reasons we chose the latter. After years of production use it has proven, at least in my opinion, to be the best decision given the state of affairs in NumPy and Python in general. The special value NaN (Not-A-Number) is used everywhere as the NA value, and there are API functions `isnull` and `notnull` which can be used across the dtypes to detect NA values.

However, it comes with it a couple of trade-offs which I most certainly have not ignored.

### 17.1.2 Support for integer NA

In the absence of high performance NA support being built into NumPy from the ground up, the primary casualty is the ability to represent NAs in integer arrays. For example:

```
In [414]: s = Series([1, 2, 3, 4, 5], index=list('abcde'))
```

```
In [415]: s
```

```
Out[415]:
```

```
a    1
b    2
c    3
d    4
e    5
```

```
In [416]: s.dtype
```

```
Out[416]: dtype('int64')
```

```
In [417]: s2 = s.reindex(['a', 'b', 'c', 'f', 'u'])
```

```
In [418]: s2
```

```
Out[418]:
```

```
a    1
b    2
```

```
c      3
f     NaN
u     NaN
```

```
In [419]: s2.dtype
Out[419]: dtype('float64')
```

This trade-off is made largely for memory and performance reasons, and also so that the resulting Series continues to be “numeric”. One possibility is to use `dtype=object` arrays instead.

### 17.1.3 NA type promotions

When introducing NAs into an existing Series or DataFrame via `reindex` or some other means, boolean and integer types will be promoted to a different dtype in order to store the NAs. These are summarized by this table:

Typeclass	Promotion dtype for storing NAs
floating	no change
object	no change
integer	cast to float64
boolean	cast to object

While this may seem like a heavy trade-off, in practice I have found very few cases where this is an issue in practice. Some explanation for the motivation here in the next section.

### 17.1.4 Why not make NumPy like R?

Many people have suggested that NumPy should simply emulate the NA support present in the more domain-specific statistical programming language R. Part of the reason is the NumPy type hierarchy:

Typeclass	Dtypes
<code>numpy.floating</code>	<code>float16</code> , <code>float32</code> , <code>float64</code> , <code>float128</code>
<code>numpy.integer</code>	<code>int8</code> , <code>int16</code> , <code>int32</code> , <code>int64</code>
<code>numpy.unsignedinteger</code>	<code>uint8</code> , <code>uint16</code> , <code>uint32</code> , <code>uint64</code>
<code>numpy.object_</code>	<code>object_</code>
<code>numpy.bool_</code>	<code>bool_</code>
<code>numpy.character</code>	<code>string_</code> , <code>unicode_</code>

The R language, by contrast, only has a handful of built-in data types: `integer`, `numeric` (floating-point), `character`, and `boolean`. NA types are implemented by reserving special bit patterns for each type to be used as the missing value. While doing this with the full NumPy type hierarchy would be possible, it would be a more substantial trade-off (especially for the 8- and 16-bit data types) and implementation undertaking.

An alternate approach is that of using masked arrays. A masked array is an array of data with an associated boolean *mask* denoting whether each value should be considered NA or not. I am personally not in love with this approach as I feel that overall it places a fairly heavy burden on the user and the library implementer. Additionally, it exacts a fairly high performance cost when working with numerical data compared with the simple approach of using `NaN`. Thus, I have chosen the Pythonic “practicality beats purity” approach and traded integer NA capability for a much simpler approach of using a special value in float and object arrays to denote NA, and promoting integer arrays to floating when NAs must be introduced.



## 17.2 Integer indexing

Label-based indexing with integer axis labels is a thorny topic. It has been discussed heavily on mailing lists and among various members of the scientific Python community. In pandas, our general viewpoint is that labels matter more than integer locations. Therefore, with an integer axis index *only* label-based indexing is possible with the standard tools like `.ix`. The following code will generate exceptions:

```
s = Series(range(5))
s[-1]
df = DataFrame(np.random.randn(5, 4))
df
df.ix[-2:]
```

This deliberate decision was made to prevent ambiguities and subtle bugs (many users reported finding bugs when the API change was made to stop “falling back” on position-based indexing).

## 17.3 Label-based slicing conventions

### 17.3.1 Non-monotonic indexes require exact matches

### 17.3.2 Endpoints are inclusive

Compared with standard Python sequence slicing in which the slice endpoint is not inclusive, label-based slicing in pandas **is inclusive**. The primary reason for this is that it is often not possible to easily determine the “successor” or next element after a particular label in an index. For example, consider the following Series:

```
In [420]: s = Series(randn(6), index=list('abcdef'))
```

```
In [421]: s
Out[421]:
a    1.337122
b   -1.531095
c    1.331458
d   -0.571329
e   -0.026671
f   -1.085663
```

Suppose we wished to slice from `c` to `e`, using integers this would be

```
In [422]: s[2:5]
Out[422]:
c    1.331458
d   -0.571329
e   -0.026671
```

However, if you only had `c` and `e`, determining the next element in the index can be somewhat complicated. For example, the following does not work:

```
s.ix['c':'e'+1]
```

A very common use case is to limit a time series to start and end at two specific dates. To enable this, we made the design design to make label-based slicing include both endpoints:

```
In [423]: s.ix['c':'e']
Out[423]:
c    1.331458
```

```
d    -0.571329
e    -0.026671
```

This is most definitely a “practicality beats purity” sort of thing, but it is something to watch out for if you expect label-based slicing to behave exactly in the way that standard Python integer slicing works.

## 17.4 Miscellaneous indexing gotchas

### 17.4.1 Reindex versus ix gotchas

Many users will find themselves using the `ix` indexing capabilities as a concise means of selecting data from a pandas object:

```
In [424]: df = DataFrame(randn(6, 4), columns=['one', 'two', 'three', 'four'],
.....:                  index=list('abcdef'))
.....:
```

```
In [425]: df
Out[425]:
```

	one	two	three	four
a	-1.114738	-0.058216	-0.486768	1.685148
b	0.112572	-1.495309	0.898435	-0.148217
c	-1.596070	0.159653	0.262136	0.036220
d	0.184735	-0.255069	-0.271020	1.288393
e	0.294633	-1.165787	0.846974	-0.685597
f	0.609099	-0.303961	0.625555	-0.059268

```
In [426]: df.ix[['b', 'c', 'e']]
Out[426]:
```

	one	two	three	four
b	0.112572	-1.495309	0.898435	-0.148217
c	-1.596070	0.159653	0.262136	0.036220
e	0.294633	-1.165787	0.846974	-0.685597

This is, of course, completely equivalent *in this case* to using the `reindex` method:

```
In [427]: df.reindex(['b', 'c', 'e'])
Out[427]:
```

	one	two	three	four
b	0.112572	-1.495309	0.898435	-0.148217
c	-1.596070	0.159653	0.262136	0.036220
e	0.294633	-1.165787	0.846974	-0.685597

Some might conclude that `ix` and `reindex` are 100% equivalent based on this. This is indeed true **except in the case of integer indexing**. For example, the above operation could alternately have been expressed as:

```
In [428]: df.ix[[1, 2, 4]]
Out[428]:
```

	one	two	three	four
b	0.112572	-1.495309	0.898435	-0.148217
c	-1.596070	0.159653	0.262136	0.036220
e	0.294633	-1.165787	0.846974	-0.685597

If you pass `[1, 2, 4]` to `reindex` you will get another thing entirely:

```
In [429]: df.reindex([1, 2, 4])
Out[429]:
```

	one	two	three	four
1	NaN	NaN	NaN	NaN
2	NaN	NaN	NaN	NaN
4	NaN	NaN	NaN	NaN

So it's important to remember that `reindex` is **strict label indexing only**. This can lead to some potentially surprising results in pathological cases where an index contains, say, both integers and strings:

```
In [430]: s = Series([1, 2, 3], index=['a', 0, 1])
```

```
In [431]: s
```

```
Out[431]:
```

```
a    1
0    2
1    3
```

```
In [432]: s.ix[[0, 1]]
```

```
Out[432]:
```

```
a    1
0    2
```

```
In [433]: s.reindex([0, 1])
```

```
Out[433]:
```

```
0    2
1    3
```

Because the index in this case does not contain solely integers, `ix` falls back on integer indexing. By contrast, `reindex` only looks for the values passed in the index, thus finding the integers 0 and 1. While it would be possible to insert some logic to check whether a passed sequence is all contained in the index, that logic would exact a very high cost in large data sets.

## 17.5 Timestamp limitations

### 17.5.1 Minimum and maximum timestamps

Since pandas represents timestamps in nanosecond resolution, the timespan that can be represented using a 64-bit integer is limited to approximately 584 years:

```
In [434]: begin = Timestamp(-9223285636854775809L)
```

```
In [435]: begin
```

```
Out[435]: <Timestamp: 1677-09-22 00:12:43.145224191>
```

```
In [436]: end = Timestamp(np.iinfo(np.int64).max)
```

```
In [437]: end
```

```
Out[437]: <Timestamp: 2262-04-11 23:47:16.854775807>
```

If you need to represent time series data outside the nanosecond timespan, use `PeriodIndex`:

```
In [438]: span = period_range('1215-01-01', '1381-01-01', freq='D')
```

```
In [439]: span
```

```
Out[439]:
```

```
<class 'pandas.tseries.period.PeriodIndex'>
freq: D
```

```
[01-Jan-1215, ..., 01-Jan-1381]
length: 60632
```

## 17.6 Parsing Dates from Text Files

When parsing multiple text file columns into a single date column, the new date column is prepended to the data and then *index\_col* specification is indexed off of the new set of columns rather than the original ones:

```
In [440]: print open('tmp.csv').read()
KORD,19990127, 19:00:00, 18:56:00, 0.8100
KORD,19990127, 20:00:00, 19:56:00, 0.0100
KORD,19990127, 21:00:00, 20:56:00, -0.5900
KORD,19990127, 21:00:00, 21:18:00, -0.9900
KORD,19990127, 22:00:00, 21:56:00, -0.5900
KORD,19990127, 23:00:00, 22:56:00, -0.5900

In [441]: date_spec = {'nominal': [1, 2], 'actual': [1, 3]}

In [442]: df = read_csv('tmp.csv', header=None,
.....:                  parse_dates=date_spec,
.....:                  keep_date_col=True,
.....:                  index_col=0)
.....:

# index_col=0 refers to the combined column "nominal" and not the original
# first column of 'KORD' strings
In [443]: df
Out[443]:
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 6 entries, 1999-01-27 19:00:00 to 1999-01-27 23:00:00
Data columns:
actual      6  non-null values
X.1         6  non-null values
X.2         6  non-null values
X.3         6  non-null values
X.4         6  non-null values
X.5         6  non-null values
dtypes: float64(1), int64(1), object(4)
```

## RPY2 / R INTERFACE

---

**Note:** This is all highly experimental. I would like to get more people involved with building a nice RPy2 interface for pandas

---

If your computer has R and rpy2 (> 2.2) installed (which will be left to the reader), you will be able to leverage the below functionality. On Windows, doing this is quite an ordeal at the moment, but users on Unix-like systems should find it quite easy. rpy2 evolves in time and the current interface is designed for the 2.2.x series, and we recommend to use over other series unless you are prepared to fix parts of the code. Released packages are available in PyPi, but should the latest code in the 2.2.x series be wanted it can be obtained with:

```
# if installing for the first time
hg clone http://bitbucket.org/lgautier/rpy2

cd rpy2
hg pull
hg update version_2.2.x
sudo python setup.py install
```

---

**Note:** To use R packages with this interface, you will need to install them inside R yourself. At the moment it cannot install them for you.

---

Once you have done installed R and rpy2, you should be able to import `pandas.rpy.common` without a hitch.

### 18.1 Transferring R data sets into Python

The `load_data` function retrieves an R data set and converts it to the appropriate pandas object (most likely a DataFrame):

```
In [982]: import pandas.rpy.common as com
```

```
In [983]: infert = com.load_data('infert')
```

```
In [984]: infert.head()
```

```
Out[984]:
```

	education	age	parity	induced	case	spontaneous	stratum	pooled.stratum
1	0-5yrs	26	6	1	1	2	1	3
2	0-5yrs	42	1	1	1	0	2	1
3	0-5yrs	39	6	2	1	0	3	4

4	0-5yrs	34	4	2	1	0	4	2
5	6-11yrs	35	3	1	1	1	5	32

## 18.2 Converting DataFrames into R objects

New in version 0.8. Starting from pandas 0.8, there is **experimental** support to convert DataFrames into the equivalent R object (that is, **data.frame**):

```
In [985]: from pandas import DataFrame

In [986]: df = DataFrame({'A': [1, 2, 3], 'B': [4, 5, 6], 'C': [7, 8, 9]},
.....:                  index=["one", "two", "three"])
.....:

In [987]: r_dataframe = com.convert_to_r_dataframe(df)

In [988]: print type(r_dataframe)
<class 'rpy2.robjjects.vectors.DataFrame'>

In [989]: print r_dataframe
      A B C
one   1 4 7
two   2 5 8
three 3 6 9
```

The DataFrame's index is stored as the `rownames` attribute of the `data.frame` instance.

You can also use `convert_to_r_matrix` to obtain a `Matrix` instance, but bear in mind that it will only work with homogeneously-typed DataFrames (as R matrices bear no information on the data type):

```
In [990]: r_matrix = com.convert_to_r_matrix(df)

In [991]: print type(r_matrix)
<class 'rpy2.robjjects.vectors.Matrix'>

In [992]: print r_matrix
      A B C
one   1 4 7
two   2 5 8
three 3 6 9
```

## 18.3 Calling R functions with pandas objects

## 18.4 High-level interface to R estimators

# RELATED PYTHON LIBRARIES

## 19.1 `la` (larry)

Keith Goodman's excellent [labeled array package](#) is very similar to pandas in many regards, though with some key differences. The main philosophical design difference is to be a wrapper around a single NumPy `ndarray` object while adding axis labeling and label-based operations and indexing. Because of this, creating a size-mutable object with heterogeneous columns (e.g. `DataFrame`) is not possible with the `la` package.

- Provide a single n-dimensional object with labeled axes with functionally analogous data alignment semantics to pandas objects
- Advanced / label-based indexing similar to that provided in pandas but setting is not supported
- Stays much closer to NumPy arrays than pandas— `larry` objects must be homogeneously typed
- `GroupBy` support is relatively limited, but a few functions are available: `group_mean`, `group_median`, and `group_ranking`
- It has a collection of analytical functions suited to quantitative portfolio construction for financial applications
- It has a collection of moving window statistics implemented in [Bottleneck](#)

## 19.2 `scikits.statsmodels`

The main [statistics and econometrics library](#) for Python. pandas has become a dependency of this library.

## 19.3 `scikits.timeseries`

`scikits.timeseries` provides a data structure for fixed frequency time series data based on the `numpy.MaskedArray` class. For time series data, it provides some of the same functionality to the pandas `Series` class. It has many more functions for time series-specific manipulation. Also, it has support for many more frequencies, though less customizable by the user (so 5-minutely data is easier to do with pandas for example).

We are aiming to merge these libraries together in the near future.





# COMPARISON WITH R / R LIBRARIES

Since pandas aims to provide a lot of the data manipulation and analysis functionality that people use R for, this page was started to provide a more detailed look at the R language and it's many 3rd party libraries as they relate to pandas. In offering comparisons with R and CRAN libraries, we care about the following things:

- **Functionality / flexibility:** what can / cannot be done with each tool
- **Performance:** how fast are operations. Hard numbers / benchmarks are preferable
- **Ease-of-use:** is one tool easier or harder to use (you may have to be the judge of this given side-by-side code comparisons)

As I do not have an encyclopedic knowledge of R packages, feel free to suggest additional CRAN packages to add to this list. This is also here to offer a big of a translation guide for users of these R packages.

## 20.1 data.frame

## 20.2 zoo

## 20.3 xts

## 20.4 plyr

## 20.5 reshape / reshape2



# API REFERENCE

## 21.1 General functions

### 21.1.1 Data manipulations

---

`pivot_table(data[, values, rows, cols, ...])` Create a spreadsheet-style pivot table as a DataFrame. The levels in the

---

#### **pandas.tools.pivot.pivot\_table**

`pandas.tools.pivot.pivot_table` (*data*, *values=None*, *rows=None*, *cols=None*, *aggfunc='mean'*,  
*fill\_value=None*, *margins=False*)

Create a spreadsheet-style pivot table as a DataFrame. The levels in the pivot table will be stored in MultiIndex objects (hierarchical indexes) on the index and columns of the result DataFrame

**Parameters** **data** : DataFrame

**values** : column to aggregate, optional

**rows** : list of column names or arrays to group on

Keys to group on the x-axis of the pivot table

**cols** : list of column names or arrays to group on

Keys to group on the y-axis of the pivot table

**aggfunc** : function, default `numpy.mean`, or list of functions

If list of functions passed, the resulting pivot table will have hierarchical columns whose top level are the function names (inferred from the function objects themselves)

**fill\_value** : scalar, default `None`

Value to replace missing values with

**margins** : boolean, default `False`

Add all row / columns (e.g. for subtotal / grand totals)

**Returns** **table** : DataFrame

## Examples

```
>>> df
   A  B  C  D
0  foo one small 1
1  foo one large 2
2  foo one large 2
3  foo two small 3
4  foo two small 3
5  bar one large 4
6  bar one small 5
7  bar two small 6
8  bar two large 7

>>> table = pivot_table(df, values='D', rows=['A', 'B'],
...                       cols=['C'], aggfunc=np.sum)
>>> table
      small large
foo one    1    4
   two    6   NaN
bar one    5    4
   two    6    7
```

---

<code>merge(left, right[, how, on, left_on, ...])</code>	Merge DataFrame objects by performing a database-style join operation by
--	--

---

<code>concat(objs[, axis, join, join_axes, ...])</code>	Concatenate pandas objects along a particular axis with optional set logic along the other a
---	--

---

## pandas.tools.merge.merge

`pandas.tools.merge.merge` (*left*, *right*, *how*='inner', *on*=None, *left\_on*=None, *right\_on*=None, *left\_index*=False, *right\_index*=False, *sort*=True, *suffixes*=('\_x', '\_y'), *copy*=True)

Merge DataFrame objects by performing a database-style join operation by columns or indexes.

If joining columns on columns, the DataFrame indexes *will be ignored*. Otherwise if joining indexes on indexes or indexes on a column or columns, the index will be passed on.

**Parameters** **left** : DataFrame

**right** : DataFrame

**how** : {'left', 'right', 'outer', 'inner'}, default 'inner'

- left: use only keys from left frame (SQL: left outer join)
- right: use only keys from right frame (SQL: right outer join)
- outer: use union of keys from both frames (SQL: full outer join)
- inner: use intersection of keys from both frames (SQL: inner join)

**on** : label or list

Field names to join on. Must be found in both DataFrames.

**left\_on** : label or list, or array-like

Field names to join on in left DataFrame. Can be a vector or list of vectors of the length of the DataFrame to use a particular vector as the join key instead of columns

**right\_on** : label or list, or array-like

Field names to join on in right DataFrame or vector/list of vectors per left\_on docs

**left\_index** : boolean, default True

Use the index from the left DataFrame as the join key(s). If it is a MultiIndex, the number of keys in the other DataFrame (either the index or a number of columns) must match the number of levels

**right\_index** : boolean, default True

Use the index from the right DataFrame as the join key. Same caveats as left\_index

**sort** : boolean, default True

Sort the join keys lexicographically in the result DataFrame

**suffixes** : 2-length sequence (tuple, list, ...)

Suffix to apply to overlapping column names in the left and right side, respectively

**copy** : boolean, default True

If False, do not copy data unnecessarily

**Returns** **merged** : DataFrame

### Examples

```
>>> A          >>> B
   lkey value    rkey value
0   foo    1      0   foo    5
1   bar    2      1   bar    6
2   baz    3      2   qux    7
3   foo    4      3   bar    8

>>> merge(A, B, left_on='lkey', right_on='rkey', how='outer')
   lkey  value_x  rkey  value_y
0   bar      2    bar      6
1   bar      2    bar      8
2   baz      3   NaN     NaN
3   foo      1    foo      5
4   foo      4    foo      5
5  NaN     NaN   qux      7
```

### pandas.tools.merge.concat

`pandas.tools.merge.concat` (*objs*, *axis*=0, *join*='outer', *join\_axes*=None, *ignore\_index*=False, *keys*=None, *levels*=None, *names*=None, *verify\_integrity*=False)

Concatenate pandas objects along a particular axis with optional set logic along the other axes. Can also add a layer of hierarchical indexing on the concatenation axis, which may be useful if the labels are the same (or overlapping) on the passed axis number

**Parameters** **objs** : list or dict of Series, DataFrame, or Panel objects

If a dict is passed, the sorted keys will be used as the *keys* argument, unless it is passed, in which case the values will be selected (see below). Any None objects will be dropped silently unless they are all None in which case an Exception will be raised

**axis** : {0, 1, ...}, default 0

The axis to concatenate along

**join** : { 'inner', 'outer' }, default 'outer'

How to handle indexes on other axis(es)

**join\_axes** : list of Index objects

Specific indexes to use for the other  $n - 1$  axes instead of performing inner/outer set logic

**verify\_integrity** : boolean, default False

Check whether the new concatenated axis contains duplicates. This can be very expensive relative to the actual data concatenation

**keys** : sequence, default None

If multiple levels passed, should contain tuples. Construct hierarchical index using the passed keys as the outermost level

**levels** : list of sequences, default None

Specific levels (unique values) to use for constructing a MultiIndex. Otherwise they will be inferred from the keys

**names** : list, default None

Names for the levels in the resulting hierarchical index

**ignore\_index** : boolean, default False

If True, do not use the index values on the concatenation axis. The resulting axis will be labeled 0, ...,  $n - 1$ . This is useful if you are concatenating objects where the concatenation axis does not have meaningful indexing information.

**Returns** **concatenated** : type of objects

### Notes

The keys, levels, and names arguments are all optional

## 21.1.2 Pickling

<code>load(path)</code>	Load pickled pandas object (or any other pickled object) from the specified
<code>save(obj, path)</code>	Pickle (serialize) object to input file path

### pandas.core.common.load

`pandas.core.common.load(path)`

Load pickled pandas object (or any other pickled object) from the specified file path

**Parameters** **path** : string

File path

**Returns** **unpickled** : type of object stored in file

**pandas.core.common.save**

`pandas.core.common.save(obj, path)`  
 Pickle (serialize) object to input file path

**Parameters** **obj** : any object

**path** : string

File path

**21.1.3 File IO**

<code>read_table(filepath_or_buffer[, sep, ...])</code>	Read general delimited file into DataFrame
<code>read_csv(filepath_or_buffer[, sep, dialect, ...])</code>	Read CSV (comma-separated) file into DataFrame
<code>ExcelFile.parse(sheetname[, header, ...])</code>	Read Excel table into DataFrame

**pandas.io.parsers.read\_table**

`pandas.io.parsers.read_table(filepath_or_buffer, sep='\t', dialect=None, header=0, index_col=None, names=None, skiprows=None, na_values=None, thousands=None, comment=None, parse_dates=False, keep_date_col=False, dayfirst=False, date_parser=None, nrows=None, iterator=False, chunksize=None, skip_footer=0, converters=None, verbose=False, delimiter=None, encoding=None, squeeze=False)`

Read general delimited file into DataFrame

Also supports optionally iterating or breaking of the file into chunks.

**Parameters** **filepath\_or\_buffer** : string or file handle / StringIO. The string could be

a URL. Valid URL schemes include http, ftp, and file. For file URLs, a host is expected. For instance, a local file could be file://localhost/path/to/table.csv

**sep** : string, default t (tab-stop)

Delimiter to use. Regular expressions are accepted.

**dialect** : string or csv.Dialect instance, default None

If None defaults to Excel dialect. Ignored if sep longer than 1 char See csv.Dialect documentation for more details

**header** : int, default 0

Row to use for the column labels of the parsed DataFrame

**skiprows** : list-like or integer

Row numbers to skip (0-indexed) or number of rows to skip (int)

**index\_col** : int or sequence, default None

Column to use as the row labels of the DataFrame. If a sequence is given, a MultiIndex is used.

**names** : array-like

List of column names

**na\_values** : list-like or dict, default None

Additional strings to recognize as NA/NaN. If dict passed, specific per-column NA values

**parse\_dates** : boolean, list of ints or names, list of lists, or dict

True -> try parsing all columns [1, 2, 3] -> try parsing columns 1, 2, 3 each as a separate date column [[1, 3]] -> combine columns 1 and 3 and parse as a single date column {'foo' : [1, 3]} -> parse columns 1, 3 as date and call result 'foo'

**keep\_date\_col** : boolean, default False

If True and parse\_dates specifies combining multiple columns then keep the original columns.

**date\_parser** : function

Function to use for converting dates to strings. Defaults to dateutil.parser

**dayfirst** : boolean, default False

DD/MM format dates, international and European format

**thousands** : str, default None

Thousands separator

**comment** : str, default None

Indicates remainder of line should not be parsed Does not support line commenting (will return empty line)

**nrows** : int, default None

Number of rows of file to read. Useful for reading pieces of large files

**iterator** : boolean, default False

Return TextParser object

**chunksize** : int, default None

Return TextParser object for iteration

**skip\_footer** : int, default 0

Number of line at bottom of file to skip

**converters** : dict. optional

Dict of functions for converting values in certain columns. Keys can either be integers or column labels

**verbose** : boolean, default False

Indicate number of NA values placed in non-numeric columns

**delimiter** : string, default None

Alternative argument name for sep. Regular expressions are accepted.

**encoding** : string, default None

Encoding to use for UTF when reading/writing (ex. 'utf-8')

**squeeze** : boolean, default False

If the parsed data only contains one column then return a Series



**Returns** **result** : DataFrame or TextParser

## pandas.io.parsers.read\_csv

```
pandas.io.parsers.read_csv(filepath_or_buffer, sep=',', dialect=None, header=0, index_col=None, names=None, skiprows=None, na_values=None, thousands=None, comment=None, parse_dates=False, keep_date_col=False, dayfirst=False, date_parser=None, nrows=None, iterator=False, chunksize=None, skip_footer=0, converters=None, verbose=False, delimiter=None, encoding=None, squeeze=False)
```

Read CSV (comma-separated) file into DataFrame

Also supports optionally iterating or breaking of the file into chunks.

**Parameters** **filepath\_or\_buffer** : string or file handle / StringIO. The string could be

a URL. Valid URL schemes include http, ftp, and file. For file URLs, a host is expected. For instance, a local file could be file ://localhost/path/to/table.csv

**sep** : string, default ','

Delimiter to use. If sep is None, will try to automatically determine this. Regular expressions are accepted.

**dialect** : string or csv.Dialect instance, default None

If None defaults to Excel dialect. Ignored if sep longer than 1 char See csv.Dialect documentation for more details

**header** : int, default 0

Row to use for the column labels of the parsed DataFrame

**skiprows** : list-like or integer

Row numbers to skip (0-indexed) or number of rows to skip (int)

**index\_col** : int or sequence, default None

Column to use as the row labels of the DataFrame. If a sequence is given, a MultiIndex is used.

**names** : array-like

List of column names

**na\_values** : list-like or dict, default None

Additional strings to recognize as NA/NaN. If dict passed, specific per-column NA values

**parse\_dates** : boolean, list of ints or names, list of lists, or dict

True -> try parsing all columns [1, 2, 3] -> try parsing columns 1, 2, 3 each as a separate date column [[1, 3]] -> combine columns 1 and 3 and parse as a single date column {'foo' : [1, 3]} -> parse columns 1, 3 as date and call result 'foo'

**keep\_date\_col** : boolean, default False

If True and parse\_dates specifies combining multiple columns then keep the original columns.

**date\_parser** : function

Function to use for converting dates to strings. Defaults to dateutil.parser

**dayfirst** : boolean, default False

DD/MM format dates, international and European format

**thousands** : str, default None

Thousands separator

**comment** : str, default None

Indicates remainder of line should not be parsed Does not support line commenting  
(will return empty line)

**nrows** : int, default None

Number of rows of file to read. Useful for reading pieces of large files

**iterator** : boolean, default False

Return TextParser object

**chunksize** : int, default None

Return TextParser object for iteration

**skip\_footer** : int, default 0

Number of line at bottom of file to skip

**converters** : dict. optional

Dict of functions for converting values in certain columns. Keys can either be integers  
or column labels

**verbose** : boolean, default False

Indicate number of NA values placed in non-numeric columns

**delimiter** : string, default None

Alternative argument name for sep. Regular expressions are accepted.

**encoding** : string, default None

Encoding to use for UTF when reading/writing (ex. 'utf-8')

**squeeze** : boolean, default False

If the parsed data only contains one column then return a Series

**Returns** **result** : DataFrame or TextParser

### pandas.io.parsers.ExcelFile.parse

ExcelFile.**parse**(*sheetname*, *header=0*, *skiprows=None*, *index\_col=None*, *parse\_cols=None*,  
*parse\_dates=False*, *date\_parser=None*, *na\_values=None*, *thousands=None*, *chunk-*  
*size=None*)

Read Excel table into DataFrame

**Parameters** **sheetname** : string

Name of Excel sheet

**header** : int, default 0

Row to use for the column labels of the parsed DataFrame

**skiprows** : list-like

Row numbers to skip (0-indexed)

**index\_col** : int, default None

Column to use as the row labels of the DataFrame. Pass None if there is no such column

**parse\_cols** : int or list, default None

If None then parse all columns, If int then indicates last column to be parsed If list of ints then indicates list of column numbers to be parsed

**na\_values** : list-like, default None

List of additional strings to recognize as NA/NaN

**Returns** **parsed** : DataFrame

#### 21.1.4 HDFStore: PyTables (HDF5)

<code>HDFStore.put(key, value[, table, append, ...])</code>	Store object in HDFStore
<code>HDFStore.get(key)</code>	Retrieve pandas object stored in file

##### pandas.io.pytables.HDFStore.put

`HDFStore.put` (*key, value, table=False, append=False, compression=None*)

Store object in HDFStore

**Parameters** **key** : object

**value** : {Series, DataFrame, Panel}

**table** : boolean, default False

Write as a PyTables Table structure which may perform worse but allow more flexible operations like searching / selecting subsets of the data

**append** : boolean, default False

For table data structures, append the input data to the existing table

**compression** : {None, 'blosc', 'lzo', 'zlib'}, default None

Use a compression algorithm to compress the data If None, the compression settings specified in the ctor will be used.

##### pandas.io.pytables.HDFStore.get

`HDFStore.get` (*key*)

Retrieve pandas object stored in file

**Parameters** **key** : object

**Returns** **obj** : type of object stored in file

#### 21.1.5 Standard moving window functions

<code>rolling_count(arg, window[, freq, time_rule])</code>	Rolling count of number of non-NaN observations inside provided window.
<code>rolling_sum(arg, window[, min_periods, ...])</code>	Moving sum
<code>rolling_mean(arg, window[, min_periods, ...])</code>	Moving mean
<code>rolling_median(arg, window[, min_periods, ...])</code>	O(N log(window)) implementation using skip list
<code>rolling_var(arg, window[, min_periods, ...])</code>	Unbiased moving variance
<code>rolling_std(arg, window[, min_periods, ...])</code>	Unbiased moving standard deviation
<code>rolling_corr(arg1, arg2, window[, ...])</code>	Moving sample correlation
<code>rolling_cov(arg1, arg2, window[, ...])</code>	Unbiased moving covariance
<code>rolling_skew(arg, window[, min_periods, ...])</code>	Unbiased moving skewness
<code>rolling_kurt(arg, window[, min_periods, ...])</code>	Unbiased moving kurtosis
<code>rolling_apply(arg, window, func[, ...])</code>	Generic moving function application
<code>rolling_quantile(arg, window, quantile[, ...])</code>	Moving quantile

### pandas.stats.moments.rolling\_count

`pandas.stats.moments.rolling_count` (*arg, window, freq=None, time\_rule=None*)  
 Rolling count of number of non-NaN observations inside provided window.

**Parameters** **arg** : DataFrame or numpy ndarray-like  
**window** : Number of observations used for calculating statistic  
**freq** : None or string alias / date offset object, default=None  
 Frequency to conform to before computing statistic

**Returns** **rolling\_count** : type of caller

### pandas.stats.moments.rolling\_sum

`pandas.stats.moments.rolling_sum` (*arg, window, min\_periods=None, freq=None, time\_rule=None, \*\*kwargs*)

Moving sum

**Parameters** **arg** : Series, DataFrame  
**window** : Number of observations used for calculating statistic  
**min\_periods** : int  
 Minimum number of observations in window required to have a value  
**freq** : None or string alias / date offset object, default=None  
 Frequency to conform to before computing statistic

**Returns** **y** : type of input argument

### pandas.stats.moments.rolling\_mean

`pandas.stats.moments.rolling_mean` (*arg, window, min\_periods=None, freq=None, time\_rule=None, \*\*kwargs*)

Moving mean

**Parameters** **arg** : Series, DataFrame  
**window** : Number of observations used for calculating statistic  
**min\_periods** : int

Minimum number of observations in window required to have a value

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**Returns** y : type of input argument

### **pandas.stats.moments.rolling\_median**

`pandas.stats.moments.rolling_median` (*arg*, *window*, *min\_periods=None*, *freq=None*,  
*time\_rule=None*, *\*\*kwargs*)

O(N log(window)) implementation using skip list

Moving median

**Parameters** **arg** : Series, DataFrame

**window** : Number of observations used for calculating statistic

**min\_periods** : int

Minimum number of observations in window required to have a value

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**Returns** y : type of input argument

### **pandas.stats.moments.rolling\_var**

`pandas.stats.moments.rolling_var` (*arg*, *window*, *min\_periods=None*, *freq=None*,  
*time\_rule=None*, *\*\*kwargs*)

Unbiased moving variance

**Parameters** **arg** : Series, DataFrame

**window** : Number of observations used for calculating statistic

**min\_periods** : int

Minimum number of observations in window required to have a value

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**Returns** y : type of input argument

### **pandas.stats.moments.rolling\_std**

`pandas.stats.moments.rolling_std` (*arg*, *window*, *min\_periods=None*, *freq=None*,  
*time\_rule=None*, *\*\*kwargs*)

Unbiased moving standard deviation

**Parameters** **arg** : Series, DataFrame

**window** : Number of observations used for calculating statistic

**min\_periods** : int

Minimum number of observations in window required to have a value

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**Returns** **y** : type of input argument

### **pandas.stats.moments.rolling\_corr**

`pandas.stats.moments.rolling_corr` (*arg1, arg2, window, min\_periods=None, time\_rule=None*)

Moving sample correlation

**Parameters** **arg1** : Series, DataFrame, or ndarray

**arg2** : Series, DataFrame, or ndarray

**window** : Number of observations used for calculating statistic

**min\_periods** : int

Minimum number of observations in window required to have a value

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**Returns** **y** : type depends on inputs

DataFrame / DataFrame -> DataFrame (matches on columns) DataFrame / Series ->

Computes result for each column Series / Series -> Series

### **pandas.stats.moments.rolling\_cov**

`pandas.stats.moments.rolling_cov` (*arg1, arg2, window, min\_periods=None, time\_rule=None*)

Unbiased moving covariance

**Parameters** **arg1** : Series, DataFrame, or ndarray

**arg2** : Series, DataFrame, or ndarray

**window** : Number of observations used for calculating statistic

**min\_periods** : int

Minimum number of observations in window required to have a value

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**Returns** **y** : type depends on inputs

DataFrame / DataFrame -> DataFrame (matches on columns) DataFrame / Series ->

Computes result for each column Series / Series -> Series

### **pandas.stats.moments.rolling\_skew**

`pandas.stats.moments.rolling_skew` (*arg, window, min\_periods=None, freq=None, time\_rule=None, \*\*kwargs*)

Unbiased moving skewness

**Parameters** **arg** : Series, DataFrame

**window** : Number of observations used for calculating statistic

**min\_periods** : int

Minimum number of observations in window required to have a value

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**Returns** **y** : type of input argument

### **pandas.stats.moments.rolling\_kurt**

`pandas.stats.moments.rolling_kurt` (*arg, window, min\_periods=None, freq=None, time\_rule=None, \*\*kwargs*)

Unbiased moving kurtosis

**Parameters** **arg** : Series, DataFrame

**window** : Number of observations used for calculating statistic

**min\_periods** : int

Minimum number of observations in window required to have a value

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**Returns** **y** : type of input argument

### **pandas.stats.moments.rolling\_apply**

`pandas.stats.moments.rolling_apply` (*arg, window, func, min\_periods=None, freq=None, time\_rule=None*)

Generic moving function application

**Parameters** **arg** : Series, DataFrame

**window** : Number of observations used for calculating statistic

**func** : function

Must produce a single value from an ndarray input

**min\_periods** : int

Minimum number of observations in window required to have a value

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**Returns** **y** : type of input argument

### **pandas.stats.moments.rolling\_quantile**

`pandas.stats.moments.rolling_quantile` (*arg, window, quantile, min\_periods=None, freq=None, time\_rule=None*)

Moving quantile

**Parameters** **arg** : Series, DataFrame

**window** : Number of observations used for calculating statistic

**quantile** :  $0 \leq \text{quantile} \leq 1$

**min\_periods** : int

Minimum number of observations in window required to have a value

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**Returns** **y** : type of input argument

## 21.1.6 Exponentially-weighted moving window functions

<code>ewma(arg[, com, span, min_periods, freq, ...])</code>	Exponentially-weighted moving average
<code>ewmstd(arg[, com, span, min_periods, bias, ...])</code>	Exponentially-weighted moving std
<code>ewmvar(arg[, com, span, min_periods, bias, ...])</code>	Exponentially-weighted moving variance
<code>ewmcorr(arg1, arg2[, com, span, ...])</code>	Exponentially-weighted moving correlation
<code>ewmcov(arg1, arg2[, com, span, min_periods, ...])</code>	Exponentially-weighted moving covariance

### **pandas.stats.moments.ewma**

`pandas.stats.moments.ewma(arg, com=None, span=None, min_periods=0, freq=None, time_rule=None, adjust=True)`

Exponentially-weighted moving average

**Parameters** **arg** : Series, DataFrame

**com** : float, optional

Center of mass:  $\alpha = \text{com} / (1 + \text{com})$ ,

**span** : float, optional

Specify decay in terms of span,  $\alpha = 2 / (\text{span} + 1)$

**min\_periods** : int, default 0

Number of observations in sample to require (only affects beginning)

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**adjust** : boolean, default True

Divide by decaying adjustment factor in beginning periods to account for imbalance in relative weightings (viewing EWMA as a moving average)

**Returns** **y** : type of input argument

### **Notes**

Either center of mass or span must be specified

EWMA is sometimes specified using a “span” parameter  $s$ , we have have that the decay parameter  $\alpha$  is related to the span as  $\alpha = 1 - 2/(s + 1) = c/(1 + c)$



where  $c$  is the center of mass. Given a span, the associated center of mass is  $c = (s - 1)/2$

So a “20-day EWMA” would have center 9.5.

### pandas.stats.moments.ewmstd

`pandas.stats.moments.ewmstd` (*arg*, *com=None*, *span=None*, *min\_periods=0*, *bias=False*,  
*time\_rule=None*)  
Exponentially-weighted moving std

**Parameters** *arg* : Series, DataFrame

**com** : float, optional

Center of mass:  $\alpha = \text{com} / (1 + \text{com})$ ,

**span** : float, optional

Specify decay in terms of span,  $\alpha = 2 / (\text{span} + 1)$

**min\_periods** : int, default 0

Number of observations in sample to require (only affects beginning)

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**adjust** : boolean, default True

Divide by decaying adjustment factor in beginning periods to account for imbalance in relative weightings (viewing EWMA as a moving average)

**bias** : boolean, default False

Use a standard estimation bias correction

**Returns** *y* : type of input argument

### Notes

Either center of mass or span must be specified

EWMA is sometimes specified using a “span” parameter  $s$ , we have have that the decay parameter  $\alpha$  is related to the span as  $\alpha = 1 - 2/(s + 1) = c/(1 + c)$

where  $c$  is the center of mass. Given a span, the associated center of mass is  $c = (s - 1)/2$

So a “20-day EWMA” would have center 9.5.

### pandas.stats.moments.ewmvar

`pandas.stats.moments.ewmvar` (*arg*, *com=None*, *span=None*, *min\_periods=0*, *bias=False*,  
*freq=None*, *time\_rule=None*)  
Exponentially-weighted moving variance

**Parameters** *arg* : Series, DataFrame

**com** : float, optional

Center of mass:  $\alpha = \text{com} / (1 + \text{com})$ ,

**span** : float, optional

Specify decay in terms of span,  $\alpha = 2 / (\text{span} + 1)$

**min\_periods** : int, default 0

Number of observations in sample to require (only affects beginning)

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**adjust** : boolean, default True

Divide by decaying adjustment factor in beginning periods to account for imbalance in relative weightings (viewing EWMA as a moving average)

**bias** : boolean, default False

Use a standard estimation bias correction

**Returns** y : type of input argument

### Notes

Either center of mass or span must be specified

EWMA is sometimes specified using a “span” parameter  $s$ , we have have that the decay parameter  $\alpha$  is related to the span as  $\alpha = 1 - 2/(s + 1) = c/(1 + c)$

where  $c$  is the center of mass. Given a span, the associated center of mass is  $c = (s - 1)/2$

So a “20-day EWMA” would have center 9.5.

### pandas.stats.moments.ewmcorr

`pandas.stats.moments.ewmcorr` (*arg1*, *arg2*, *com=None*, *span=None*, *min\_periods=0*, *freq=None*,  
*time\_rule=None*)

Exponentially-weighted moving correlation

**Parameters** **arg1** : Series, DataFrame, or ndarray

**arg2** : Series, DataFrame, or ndarray

**com** : float. optional

Center of mass:  $\alpha = \text{com} / (1 + \text{com})$ ,

**span** : float, optional

Specify decay in terms of span,  $\alpha = 2 / (\text{span} + 1)$

**min\_periods** : int, default 0

Number of observations in sample to require (only affects beginning)

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**adjust** : boolean, default True

Divide by decaying adjustment factor in beginning periods to account for imbalance in relative weightings (viewing EWMA as a moving average)

**Returns** y : type of input argument

## Notes

Either center of mass or span must be specified

EWMA is sometimes specified using a “span” parameter  $s$ , we have have that the decay parameter  $\alpha$  is related to the span as  $\alpha = 1 - 2/(s + 1) = c/(1 + c)$

where  $c$  is the center of mass. Given a span, the associated center of mass is  $c = (s - 1)/2$

So a “20-day EWMA” would have center 9.5.

## pandas.stats.moments.ewmcov

`pandas.stats.moments.ewmcov` (*arg1*, *arg2*, *com=None*, *span=None*, *min\_periods=0*, *bias=False*,  
*freq=None*, *time\_rule=None*)  
Exponentially-weighted moving covariance

**Parameters** **arg1** : Series, DataFrame, or ndarray

**arg2** : Series, DataFrame, or ndarray

**com** : float, optional

Center of mass:  $\alpha = \text{com} / (1 + \text{com})$ ,

**span** : float, optional

Specify decay in terms of span,  $\alpha = 2 / (\text{span} + 1)$

**min\_periods** : int, default 0

Number of observations in sample to require (only affects beginning)

**freq** : None or string alias / date offset object, default=None

Frequency to conform to before computing statistic

**adjust** : boolean, default True

Divide by decaying adjustment factor in beginning periods to account for imbalance in relative weightings (viewing EWMA as a moving average)

**Returns** **y** : type of input argument

## Notes

Either center of mass or span must be specified

EWMA is sometimes specified using a “span” parameter  $s$ , we have have that the decay parameter  $\alpha$  is related to the span as  $\alpha = 1 - 2/(s + 1) = c/(1 + c)$

where  $c$  is the center of mass. Given a span, the associated center of mass is  $c = (s - 1)/2$

So a “20-day EWMA” would have center 9.5.

## 21.2 Series

### 21.2.1 Attributes and underlying data

#### Axes

- **index:** axis labels

<code>Series.values</code>	Return Series as ndarray
<code>Series.dtype</code>	Data-type of the array's elements.
<code>Series.isnull(obj)</code>	Replacement for <code>numpy.isnan</code> / <code>-numpy.isfinite</code> which is suitable for use on object arrays.
<code>Series.notnull(obj)</code>	Replacement for <code>numpy.isfinite</code> / <code>-numpy.isnan</code> which is suitable for use on object arrays.

## pandas.Series.values

### `Series.values`

Return Series as ndarray

**Returns** `arr` : `numpy.ndarray`

## pandas.Series.dtype

### `Series.dtype`

Data-type of the array's elements.

**Parameters** `None` :

**Returns** `d` : `numpy dtype object`

**See Also:**

`numpy.dtype`

### Examples

```
>>> x
array([[0, 1],
       [2, 3]])
>>> x.dtype
dtype('int32')
>>> type(x.dtype)
<type 'numpy.dtype'>
```

## pandas.Series.isnull

### `Series.isnull(obj)`

Replacement for `numpy.isnan` / `-numpy.isfinite` which is suitable for use on object arrays.

**Parameters** `arr`: ndarray or object value :

**Returns** `boolean ndarray or boolean` :

## pandas.Series.notnull

### `Series.notnull(obj)`

Replacement for `numpy.isfinite` / `-numpy.isnan` which is suitable for use on object arrays.

**Parameters** `arr`: ndarray or object value :

**Returns** `boolean ndarray or boolean` :

## 21.2.2 Conversion / Constructors

<code>Series.__init__([data, index, dtype, name, copy])</code>	One-dimensional ndarray with axis labels (including time series).
<code>Series.astype(dtype)</code>	See <code>numpy.ndarray.astype</code>
<code>Series.copy([order])</code>	Return new Series with copy of underlying values

### pandas.Series.\_\_init\_\_

`Series.__init__(data=None, index=None, dtype=None, name=None, copy=False)`

One-dimensional ndarray with axis labels (including time series). Labels need not be unique but must be any hashable type. The object supports both integer- and label-based indexing and provides a host of methods for performing operations involving the index. Statistical methods from ndarray have been overridden to automatically exclude missing data (currently represented as NaN)

Operations between Series (+, -, /, \*) align values based on their associated index values– they need not be the same length. The result index will be the sorted union of the two indexes.

**Parameters** **data** : array-like, dict, or scalar value

Contains data stored in Series

**index** : array-like or Index (1d)

Values must be unique and hashable, same length as data. Index object (or other iterable of same length as data) Will default to `np.arange(len(data))` if not provided. If both a dict and index sequence are used, the index will override the keys found in the dict.

**dtype** : numpy.dtype or None

If None, dtype will be inferred copy : boolean, default False Copy input data

**copy** : boolean, default False

### pandas.Series.astype

`Series.astype(dtype)`

See `numpy.ndarray.astype`

### pandas.Series.copy

`Series.copy(order='C')`

Return new Series with copy of underlying values

**Returns** **cp** : Series

## 21.2.3 Indexing, iteration

<code>Series.get(label[, default])</code>	Returns value occupying requested label, default to specified missing value if not present.
<code>Series.ix</code>	
<code>Series.__iter__()</code>	
<code>Series.iteritems([index])</code>	Lazily iterate over (index, value) tuples

## pandas.Series.get

`Series.get (label, default=None)`

Returns value occupying requested label, default to specified missing value if not present. Analogous to dict.get

**Parameters** **label** : object

Label value looking for

**default** : object, optional

Value to return if label not in index

**Returns** **y** : scalar

## pandas.Series.ix

`Series.ix`

## pandas.Series.\_\_iter\_\_

`Series.__iter__()`

## pandas.Series.iteritems

`Series.iteritems (index=True)`

Lazily iterate over (index, value) tuples

## 21.2.4 Binary operator functions

<code>Series.add(other[, level, fill_value])</code>	Binary operator add with support to substitute a fill_value for missing data
<code>Series.div(other[, level, fill_value])</code>	Binary operator divide with support to substitute a fill_value for missing data
<code>Series.mul(other[, level, fill_value])</code>	Binary operator multiply with support to substitute a fill_value for missing data
<code>Series.sub(other[, level, fill_value])</code>	Binary operator subtract with support to substitute a fill_value for missing data
<code>Series.combine(other, func[, fill_value])</code>	Perform elementwise binary operation on two Series using given function
<code>Series.combine_first(other)</code>	Combine Series values, choosing the calling Series's values

## pandas.Series.add

`Series.add (other, level=None, fill_value=None)`

Binary operator add with support to substitute a fill\_value for missing data in one of the inputs

**Parameters** **other**: Series or scalar value :

**fill\_value** : None or float value, default None (NaN)

Fill missing (NaN) values with this value. If both Series are missing, the result will be missing

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**Returns** **result** : Series

### pandas.Series.div

`Series.div` (*other*, *level=None*, *fill\_value=None*)

Binary operator divide with support to substitute a *fill\_value* for missing data in one of the inputs

**Parameters** **other:** Series or scalar value :

**fill\_value** : None or float value, default None (NaN)

Fill missing (NaN) values with this value. If both Series are missing, the result will be missing

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**Returns** **result** : Series

### pandas.Series.mul

`Series.mul` (*other*, *level=None*, *fill\_value=None*)

Binary operator multiply with support to substitute a *fill\_value* for missing data in one of the inputs

**Parameters** **other:** Series or scalar value :

**fill\_value** : None or float value, default None (NaN)

Fill missing (NaN) values with this value. If both Series are missing, the result will be missing

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**Returns** **result** : Series

### pandas.Series.sub

`Series.sub` (*other*, *level=None*, *fill\_value=None*)

Binary operator subtract with support to substitute a *fill\_value* for missing data in one of the inputs

**Parameters** **other:** Series or scalar value :

**fill\_value** : None or float value, default None (NaN)

Fill missing (NaN) values with this value. If both Series are missing, the result will be missing

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**Returns** **result** : Series

### pandas.Series.combine

`Series.combine` (*other*, *func*, *fill\_value=nan*)

Perform elementwise binary operation on two Series using given function with optional fill value when an index is missing from one Series or the other

**Parameters** **other** : Series or scalar value

**func** : function

**fill\_value** : scalar value

**Returns** **result** : Series

### **pandas.Series.combine\_first**

**Series.combine\_first** (*other*)

Combine Series values, choosing the calling Series's values first. Result index will be the union of the two indexes

**Parameters** **other** : Series

**Returns** **y** : Series

## **21.2.5 Function application, GroupBy**

<code>Series.apply(func[, convert_dtype])</code>	Invoke function on values of Series. Can be ufunc or Python function
<code>Series.map(arg)</code>	Map values of Series using input correspondence (which can be
<code>Series.groupby([by, axis, level, as_index, ...])</code>	Group series using mapper (dict or key function, apply given function

### **pandas.Series.apply**

**Series.apply** (*func, convert\_dtype=True*)

Invoke function on values of Series. Can be ufunc or Python function expecting only single values

**Parameters** **func** : function

**convert\_dtype** : boolean, default True

Try to find better dtype for elementwise function results. If False, leave as dtype=object

**Returns** **y** : Series

**See Also:**

**Series.map** For element-wise operations

### **pandas.Series.map**

**Series.map** (*arg*)

Map values of Series using input correspondence (which can be a dict, Series, or function)

**Parameters** **arg** : function, dict, or Series

**Returns** **y** : Series

same index as caller

### **Examples**



```

>>> x
one    1
two    2
three  3

>>> y
1    foo
2    bar
3    baz

>>> x.map(y)
one    foo
two    bar
three  baz

```

## pandas.Series.groupby

`Series.groupby` (*by=None, axis=0, level=None, as\_index=True, sort=True, group\_keys=True*)

Group series using mapper (dict or key function, apply given function to group, return result as series) or by a series of columns

**Parameters** **by** : mapping function / list of functions, dict, Series, or tuple /

list of column names. Called on each element of the object index to determine the groups. If a dict or Series is passed, the Series or dict VALUES will be used to determine the groups

**axis** : int, default 0

**level** : int, level name, or sequence of such, default None

If the axis is a MultiIndex (hierarchical), group by a particular level or levels

**as\_index** : boolean, default True

For aggregated output, return object with group labels as the index. Only relevant for DataFrame input. `as_index=False` is effectively “SQL-style” grouped output

**sort** : boolean, default True

Sort group keys. Get better performance by turning this off

**group\_keys** : boolean, default True

When calling `apply`, add group keys to index to identify pieces

**Returns** GroupBy object :

## Examples

```

# DataFrame result >>> data.groupby(func, axis=0).mean()
# DataFrame result >>> data.groupby(['col1', 'col2'])['col3'].mean()
# DataFrame with hierarchical index >>> data.groupby(['col1', 'col2']).mean()

```

## 21.2.6 Computations / Descriptive Stats

<code>Series.autocorr()</code>	Lag-1 autocorrelation
<code>Series.clip([lower, upper, out])</code>	Trim values at input threshold(s)
<code>Series.clip_lower(threshold)</code>	Return copy of series with values below given value truncated
<code>Series.clip_upper(threshold)</code>	Return copy of series with values above given value truncated
<code>Series.corr(other[, method])</code>	Compute correlation two Series, excluding missing values
<code>Series.count([level])</code>	Return number of non-NA/null observations in the Series
<code>Series.cumprod([axis, dtype, out, skipna])</code>	Cumulative product of values.
<code>Series.cumsum([axis, dtype, out, skipna])</code>	Cumulative sum of values.
<code>Series.describe([percentile_width])</code>	Generate various summary statistics of Series, excluding NaN
<code>Series.diff([periods])</code>	1st discrete difference of object
<code>Series.max([axis, out, skipna, level])</code>	Return maximum of values
<code>Series.mean([axis, dtype, out, skipna, level])</code>	Return mean of values
<code>Series.median([axis, dtype, out, skipna, level])</code>	Return median of values
<code>Series.min([axis, out, skipna, level])</code>	Return minimum of values
<code>Series.prod([axis, dtype, out, skipna, level])</code>	Return product of values
<code>Series.quantile([q])</code>	Return value at the given quantile, a la scoreatpercentile in
<code>Series.skew([skipna, level])</code>	Return unbiased skewness of values
<code>Series.std([axis, dtype, out, ddof, skipna, ...])</code>	Return standard deviation of values
<code>Series.sum([axis, dtype, out, skipna, level])</code>	Return sum of values
<code>Series.var([axis, dtype, out, ddof, skipna, ...])</code>	Return variance of values
<code>Series.value_counts()</code>	Returns Series containing counts of unique values. The resulting Series

### **pandas.Series.autocorr**

`Series.autocorr()`  
Lag-1 autocorrelation  
**Returns** `autocorr` : float

### **pandas.Series.clip**

`Series.clip(lower=None, upper=None, out=None)`  
Trim values at input threshold(s)  
**Parameters** `lower` : float, default None  
`upper` : float, default None  
**Returns** `clipped` : Series

### **pandas.Series.clip\_lower**

`Series.clip_lower(threshold)`  
Return copy of series with values below given value truncated  
**Returns** `clipped` : Series  
**See Also:**  
`clip`

### pandas.Series.clip\_upper

`Series.clip_upper` (*threshold*)

Return copy of series with values above given value truncated

**Returns** `clipped` : Series

**See Also:**

`clip`

### pandas.Series.corr

`Series.corr` (*other*, *method='pearson'*)

Compute correlation two Series, excluding missing values

**Parameters** `other` : Series

`method` : { 'pearson', 'kendall', 'spearman' }

pearson : standard correlation coefficient    kendall : Kendall Tau correlation coefficient

spearman : Spearman rank correlation

**Returns** `correlation` : float

### pandas.Series.count

`Series.count` (*level=None*)

Return number of non-NA/null observations in the Series

**Parameters** `level` : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a smaller Series

**Returns** `nobs` : int or Series (if level specified)

### pandas.Series.cumprod

`Series.cumprod` (*axis=0*, *dtype=None*, *out=None*, *skipna=True*)

Cumulative product of values. Preserves locations of NaN values

Extra parameters are to preserve ndarray interface.

**Parameters** `skipna` : boolean, default True

Exclude NA/null values

**Returns** `cumprod` : Series

### pandas.Series.cumsum

`Series.cumsum` (*axis=0*, *dtype=None*, *out=None*, *skipna=True*)

Cumulative sum of values. Preserves locations of NaN values

Extra parameters are to preserve ndarray interface.

**Parameters** `skipna` : boolean, default True

Exclude NA/null values

**Returns** `cumsum` : Series

### **pandas.Series.describe**

`Series.describe` (*percentile\_width=50*)

Generate various summary statistics of Series, excluding NaN values. These include: count, mean, std, min, max, and lower%/50%/upper% percentiles

**Parameters** `percentile_width` : float, optional

width of the desired uncertainty interval, default is 50, which corresponds to lower=25, upper=75

**Returns** `desc` : Series

### **pandas.Series.diff**

`Series.diff` (*periods=1*)

1st discrete difference of object

**Parameters** `periods` : int, default 1

Periods to shift for forming difference

**Returns** `dified` : Series

### **pandas.Series.max**

`Series.max` (*axis=None, out=None, skipna=True, level=None*)

Return maximum of values NA/null values are excluded

**Parameters** `skipna` : boolean, default True

Exclude NA/null values

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a smaller Series

**Returns** `max` : float (or Series if level specified)

### **pandas.Series.mean**

`Series.mean` (*axis=0, dtype=None, out=None, skipna=True, level=None*)

Return mean of values NA/null values are excluded

**Parameters** `skipna` : boolean, default True

Exclude NA/null values

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a smaller Series

**Extra parameters are to preserve ndarrayinterface. :**

**Returns** `mean` : float (or Series if level specified)

### pandas.Series.median

`Series.median` (*axis=0, dtype=None, out=None, skipna=True, level=None*)

Return median of values NA/null values are excluded

**Parameters** `skipna` : boolean, default True

Exclude NA/null values

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a smaller Series

**Returns** `median` : float (or Series if level specified)

### pandas.Series.min

`Series.min` (*axis=None, out=None, skipna=True, level=None*)

Return minimum of values NA/null values are excluded

**Parameters** `skipna` : boolean, default True

Exclude NA/null values

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a smaller Series

**Returns** `min` : float (or Series if level specified)

### pandas.Series.prod

`Series.prod` (*axis=0, dtype=None, out=None, skipna=True, level=None*)

Return product of values NA/null values are excluded

**Parameters** `skipna` : boolean, default True

Exclude NA/null values

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a smaller Series

**Returns** `prod` : float (or Series if level specified)

### pandas.Series.quantile

`Series.quantile` (*q=0.5*)

Return value at the given quantile, a la `scoreatpercentile` in `scipy.stats`

**Parameters** `q` : quantile

$0 \leq q \leq 1$

**Returns** `quantile` : float

### pandas.Series.skew

`Series.skew(skipna=True, level=None)`

Return unbiased skewness of values NA/null values are excluded

**Parameters** `skipna` : boolean, default True

Exclude NA/null values

`level` : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a smaller Series

**Returns** `skew` : float (or Series if level specified)

### pandas.Series.std

`Series.std(axis=None, dtype=None, out=None, ddof=1, skipna=True, level=None)`

Return standard deviation of values NA/null values are excluded

**Parameters** `skipna` : boolean, default True

Exclude NA/null values

`level` : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a smaller Series

**Returns** `stdev` : float (or Series if level specified)

### pandas.Series.sum

`Series.sum(axis=0, dtype=None, out=None, skipna=True, level=None)`

Return sum of values NA/null values are excluded

**Parameters** `skipna` : boolean, default True

Exclude NA/null values

`level` : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a smaller Series

**Extra parameters are to preserve ndarrayinterface. :**

**Returns** `sum` : float (or Series if level specified)

### pandas.Series.var

`Series.var(axis=None, dtype=None, out=None, ddof=1, skipna=True, level=None)`

Return variance of values NA/null values are excluded

**Parameters** `skipna` : boolean, default True

Exclude NA/null values

`level` : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a smaller Series

**Returns** **var** : float (or Series if level specified)

## pandas.Series.value\_counts

`Series.value_counts()`

Returns Series containing counts of unique values. The resulting Series will be in descending order so that the first element is the most frequently-occurring element. Excludes NA values

**Returns** **counts** : Series

## 21.2.7 Reindexing / Selection / Label manipulation

<code>Series.align(other[, join, level, copy, ...])</code>	Align two Series object with the specified join method
<code>Series.drop(labels[, axis, level])</code>	Return new object with labels in requested axis removed
<code>Series.reindex([index, method, level, ...])</code>	Conform Series to new index with optional filling logic, placing
<code>Series.reindex_like(other[, method, limit])</code>	Reindex Series to match index of another Series, optionally with
<code>Series.rename(mapper[, inplace])</code>	Alter Series index using dict or function
<code>Series.select(crit[, axis])</code>	Return data corresponding to axis labels matching criteria
<code>Series.take(indices[, axis])</code>	Analogous to ndarray.take, return Series corresponding to requested
<code>Series.truncate([before, after, copy])</code>	Function truncate a sorted DataFrame / Series before and/or after

## pandas.Series.align

`Series.align(other, join='outer', level=None, copy=True, fill_value=None, method=None, inplace=False, limit=None)`

Align two Series object with the specified join method

**Parameters** **other** : Series

**join** : { 'outer', 'inner', 'left', 'right' }, default 'outer'

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**copy** : boolean, default True

Always return new objects. If copy=False and no reindexing is required, the same object will be returned (for better performance)

**fill\_value** : object, default None

**method** : str, default 'pad'

**limit** : int, default None

fill\_value, method, inplace, limit are passed to fillna

**Returns** **(left, right)** : (Series, Series)

Aligned Series

## pandas.Series.drop

`Series.drop(labels, axis=0, level=None)`

Return new object with labels in requested axis removed

**Parameters** **labels** : array-like

**axis** : int

**level** : int or name, default None

For MultiIndex

**Returns** **dropped** : type of caller

## pandas.Series.reindex

`Series.reindex(index=None, method=None, level=None, fill_value=nan, limit=None, copy=True)`

Conform Series to new index with optional filling logic, placing NA/NaN in locations having no value in the previous index. A new object is produced unless the new index is equivalent to the current one and copy=False

**Parameters** **index** : array-like or Index

New labels / index to conform to. Preferably an Index object to avoid duplicating data

**method** : {‘backfill’, ‘bfill’, ‘pad’, ‘ffill’, None}

Method to use for filling holes in reindexed Series pad / ffill: propagate LAST valid observation forward to next valid backfill / bfill: use NEXT valid observation to fill gap

**copy** : boolean, default True

Return a new object, even if the passed indexes are the same

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**fill\_value** : scalar, default np.NaN

Value to use for missing values. Defaults to NaN, but can be any “compatible” value

**limit** : int, default None

Maximum size gap to forward or backward fill

**Returns** **reindexed** : Series

## pandas.Series.reindex\_like

`Series.reindex_like(other, method=None, limit=None)`

Reindex Series to match index of another Series, optionally with filling logic

**Parameters** **other** : Series

**method** : string or None

See Series.reindex docstring

**limit** : int, default None

Maximum size gap to forward or backward fill

**Returns** **reindexed** : Series



## Notes

Like calling `s.reindex(other.index, method=...)`

## pandas.Series.rename

`Series.rename` (*mapper, inplace=False*)

Alter Series index using dict or function

**Parameters** `mapper` : dict-like or function

Transformation to apply to each index

**Returns** `renamed` : Series (new object)

## Notes

Function / dict values must be unique (1-to-1)

## Examples

```
>>> x
foo 1
bar 2
baz 3

>>> x.rename(str.upper)
FOO 1
BAR 2
BAZ 3

>>> x.rename({'foo' : 'a', 'bar' : 'b', 'baz' : 'c'})
a 1
b 2
c 3
```

## pandas.Series.select

`Series.select` (*crit, axis=0*)

Return data corresponding to axis labels matching criteria

**Parameters** `crit` : function

To be called on each index (label). Should return True or False

`axis` : int

**Returns** `selection` : type of caller

## pandas.Series.take

`Series.take` (*indices, axis=0*)

Analogous to `ndarray.take`, return Series corresponding to requested indices

**Parameters** `indices` : list / array of ints

**Returns** **taken** : Series

## pandas.Series.truncate

Series.**truncate** (*before=None, after=None, copy=True*)

Function truncate a sorted DataFrame / Series before and/or after some particular dates.

**Parameters** **before** : date

Truncate before date

**after** : date

Truncate after date

**Returns** **truncated** : type of caller

## 21.2.8 Missing data handling

<code>Series.dropna()</code>	Return Series without null values
<code>Series.fillna([value, method, inplace, limit])</code>	Fill NA/NaN values using the specified method
<code>Series.interpolate([method])</code>	Interpolate missing values (after the first valid value)

## pandas.Series.dropna

Series.**dropna** ()

Return Series without null values

**Returns** **valid** : Series

## pandas.Series.fillna

Series.**fillna** (*value=None, method='pad', inplace=False, limit=None*)

Fill NA/NaN values using the specified method

**Parameters** **value** : any kind (should be same type as array)

Value to use to fill holes (e.g. 0)

**method** : { 'backfill', 'bfill', 'pad', 'ffill', None }, default 'pad'

Method to use for filling holes in reindexed Series pad / ffill: propagate last valid observation forward to next valid backfill / bfill: use NEXT valid observation to fill gap

**inplace** : boolean, default False

If True, fill the Series in place. Note: this will modify any other views on this Series, for example a column in a DataFrame. Returns a reference to the filled object, which is self if inplace=True

**limit** : int, default None

Maximum size gap to forward or backward fill

**Returns** **filled** : Series

**See Also:**

`reindex, asfreq`

**pandas.Series.interpolate**`Series.interpolate (method='linear')`

Interpolate missing values (after the first valid value)

**Parameters** `method` : {'linear', 'time', 'values'}

Interpolation method. 'time' interpolation works on daily and higher resolution data to interpolate given length of interval 'values' using the actual index numeric values

**Returns** `interpolated` : Series**21.2.9 Reshaping, sorting**

<code>Series.argsort([axis, kind, order])</code>	Overrides ndarray.argsort.
<code>Series.order([na_last, ascending, kind])</code>	Sorts Series object, by value, maintaining index-value link
<code>Series.sort([axis, kind, order])</code>	Sort values and index labels by value, in place.
<code>Series.sort_index([ascending])</code>	Sort object by labels (along an axis)
<code>Series.sortlevel([level, ascending])</code>	Sort Series with MultiIndex by chosen level. Data will be
<code>Series.unstack([level])</code>	Unstack, a.k.a.

**pandas.Series.argsort**`Series.argsort (axis=0, kind='quicksort', order=None)`

Overrides ndarray.argsort. Argsorts the value, omitting NA/null values, and places the result in the same locations as the non-NA values

**Parameters** `axis` : int (can only be zero)`kind` : {'mergesort', 'quicksort', 'heapsort'}, default 'quicksort'

Choice of sorting algorithm. See np.sort for more information. 'mergesort' is the only stable algorithm

`order` : ignored**Returns** `argsorted` : Series**pandas.Series.order**`Series.order (na_last=True, ascending=True, kind='mergesort')`

Sorts Series object, by value, maintaining index-value link

**Parameters** `na_last` : boolean (optional, default=True)

Put NaN's at beginning or end

`ascending` : boolean, default True

Sort ascending. Passing False sorts descending

`kind` : {'mergesort', 'quicksort', 'heapsort'}, default 'mergesort'

Choice of sorting algorithm. See np.sort for more information. 'mergesort' is the only stable algorithm

**Returns** `y` : Series

## pandas.Series.sort

`Series.sort` (*axis=0, kind='quicksort', order=None*)

Sort values and index labels by value, in place. For compatibility with ndarray API. No return value

**Parameters** `axis` : int (can only be zero)

`kind` : { 'mergesort', 'quicksort', 'heapsort' }, default 'quicksort'

Choice of sorting algorithm. See `np.sort` for more information. 'mergesort' is the only stable algorithm

`order` : ignored

## pandas.Series.sort\_index

`Series.sort_index` (*ascending=True*)

Sort object by labels (along an axis)

**Parameters** `ascending` : boolean, default True

Sort ascending vs. descending

**Returns** `sorted_obj` : Series

## pandas.Series.sortlevel

`Series.sortlevel` (*level=0, ascending=True*)

Sort Series with MultiIndex by chosen level. Data will be lexicographically sorted by the chosen level followed by the other levels (in order)

**Parameters** `level` : int

`ascending` : bool, default True

**Returns** `sorted` : Series

## pandas.Series.unstack

`Series.unstack` (*level=-1*)

Unstack, a.k.a. pivot, Series with MultiIndex to produce DataFrame

**Parameters** `level` : int, string, or list of these, default last level

Level(s) to unstack, can pass level name

**Returns** `unstacked` : DataFrame

## Examples

```
>>> s
one  a    1.
one  b    2.
two  a    3.
two  b    4.
```

```
>>> s.unstack(level=-1)
      a    b
one  1.  2.
two  3.  4.

>>> s.unstack(level=0)
      one  two
a    1.   2.
b    3.   4.
```

## 21.2.10 Combining / joining / merging

---

`Series.append(to_append[, verify_integrity])` Concatenate two or more Series. The indexes must not overlap

---

### pandas.Series.append

`Series.append(to_append, verify_integrity=False)`

Concatenate two or more Series. The indexes must not overlap

**Parameters** `to_append` : Series or list/tuple of Series

`verify_integrity` : boolean, default False

If True, raise Exception on creating index with duplicates

**Returns** `appended` : Series

## 21.2.11 Time series-related

<code>Series.asfreq(freq[, method, how])</code>	Convert all TimeSeries inside to specified frequency using DateOffset
<code>Series.asof(when)</code>	Return last good (non-NaN) value in TimeSeries if value is NaN for
<code>Series.shift([periods, freq])</code>	Shift the index of the Series by desired number of periods with an
<code>Series.first_valid_index()</code>	Return label for first non-NA/null value
<code>Series.last_valid_index()</code>	Return label for last non-NA/null value
<code>Series.weekday</code>	

---

### pandas.Series.asfreq

`Series.asfreq(freq, method=None, how=None)`

Convert all TimeSeries inside to specified frequency using DateOffset objects. Optionally provide fill method to pad/backfill missing values.

**Parameters** `freq` : DateOffset object, or string

`method` : { 'backfill', 'bfill', 'pad', 'ffill', None }

Method to use for filling holes in reindexed Series pad / ffill: propagate last valid observation forward to next valid backfill / bfill: use NEXT valid observation to fill methdo

`how` : { 'start', 'end' }, default end

For PeriodIndex only, see PeriodIndex.asfreq

**Returns** `converted` : type of caller

## pandas.Series.asof

`Series.asof (where)`

Return last good (non-NaN) value in TimeSeries if value is NaN for requested date.

If there is no good value, NaN is returned.

**Parameters** `where` : date or array of dates

**Returns** `value or NaN` :

### Notes

Dates are assumed to be sorted

## pandas.Series.shift

`Series.shift (periods=1, freq=None, **kws)`

Shift the index of the Series by desired number of periods with an optional time offset

**Parameters** `periods` : int

Number of periods to move, can be positive or negative

**freq** : DateOffset, timedelta, or offset alias string, optional

Increment to use from datetools module or time rule (e.g. 'EOM')

**Returns** `shifted` : Series

## pandas.Series.first\_valid\_index

`Series.first_valid_index()`

Return label for first non-NA/null value

## pandas.Series.last\_valid\_index

`Series.last_valid_index()`

Return label for last non-NA/null value

## pandas.Series.weekday

`Series.weekday`

## 21.2.12 Plotting

---

<code>Series.hist([ax, grid, xlabelsize, xrot, ...])</code>	Draw histogram of the input series using matplotlib
<code>Series.plot(series[, label, kind, ...])</code>	Plot the input series with the index on the x-axis using matplotlib

---

## pandas.Series.hist

`Series.hist` (*ax=None, grid=True, xlabelsize=None, xrot=None, ylabelsize=None, yrot=None, \*\*kws*)  
Draw histogram of the input series using matplotlib

**Parameters** **ax** : matplotlib axis object

If not passed, uses `gca()`

**grid** : boolean, default True

Whether to show axis grid lines

**xlabelsize** : int, default None

If specified changes the x-axis label size

**xrot** : float, default None

rotation of x axis labels

**ylabelsize** : int, default None

If specified changes the y-axis label size

**yrot** : float, default None

rotation of y axis labels

**kws** : keywords

To be passed to the actual plotting function

### Notes

See matplotlib documentation online for more on this

## pandas.Series.plot

`Series.plot` (*series, label=None, kind='line', use\_index=True, rot=None, xticks=None, yticks=None, xlim=None, ylim=None, ax=None, style=None, grid=None, logy=False, secondary\_y=False, \*\*kws*)

Plot the input series with the index on the x-axis using matplotlib

**Parameters** **label** : label argument to provide to plot

**kind** : { 'line', 'bar' }

**rot** : int, default 30

Rotation for tick labels

**use\_index** : boolean, default True

Plot index as axis tick labels

**ax** : matplotlib axis object

If not passed, uses `gca()`

**style** : string, default matplotlib default

matplotlib line style to use

**ax** : matplotlib axis object

If not passed, uses `gca()`

**kind** : { 'line', 'bar', 'barh' }

bar : vertical bar plot barh : horizontal bar plot

**logy** : boolean, default False

For line plots, use log scaling on y axis

**xticks** : sequence

Values to use for the xticks

**yticks** : sequence

Values to use for the yticks

**xlim** : 2-tuple/list

**ylim** : 2-tuple/list

**rot** : int, default None

Rotation for ticks

**kws** : keywords

Options to pass to matplotlib plotting method

### Notes

See matplotlib documentation online for more on this subject

## 21.2.13 Serialization / IO / Conversion

<code>Series.from_csv(path[, sep, parse_dates, ...])</code>	Read delimited file into Series
<code>Series.load(path)</code>	
<code>Series.save(path)</code>	
<code>Series.to_csv(path[, index, sep, na_rep, ...])</code>	Write Series to a comma-separated values (csv) file
<code>Series.to_dict()</code>	Convert Series to {label -> value} dict
<code>Series.to_sparse([kind, fill_value])</code>	Convert Series to SparseSeries

### pandas.Series.from\_csv

**classmethod** `Series.from_csv` (*path*, *sep*=' ', *parse\_dates*=True, *header*=None, *index\_col*=0, *encoding*=None)

Read delimited file into Series

**Parameters** **path** : string file path or file handle / StringIO

**sep** : string, default ' '

Field delimiter

**parse\_dates** : boolean, default True

Parse dates. Different default from `read_table`

**header** : int, default 0



Row to use at header (skip prior rows)

**index\_col** : int or sequence, default 0

Column to use for index. If a sequence is given, a MultiIndex is used. Different default from read\_table

**encoding** : string, optional

a string representing the encoding to use if the contents are non-ascii, for python versions prior to 3

**Returns** y : Series

### pandas.Series.load

**classmethod** Series.load(path)

### pandas.Series.save

Series.save(path)

### pandas.Series.to\_csv

Series.to\_csv(path, index=True, sep=',', na\_rep='', header=False, index\_label=None, mode='w', nan-Rep=None, encoding=None)

Write Series to a comma-separated values (csv) file

**Parameters** path : string file path or file handle / StringIO

na\_rep : string, default ''

Missing data rep'n

header : boolean, default False

Write out series name

index : boolean, default True

Write row names (index)

index\_label : string or sequence, default None

Column label for index column(s) if desired. If None is given, and header and index are True, then the index names are used. A sequence should be given if the DataFrame uses MultiIndex.

mode : Python write mode, default 'w'

sep : character, default ','

Field delimiter for the output file.

encoding : string, optional

a string representing the encoding to use if the contents are non-ascii, for python versions prior to 3

## pandas.Series.to\_dict

`Series.to_dict()`  
Convert Series to {label -> value} dict  
**Returns** `value_dict` : dict

## pandas.Series.to\_sparse

`Series.to_sparse(kind='block', fill_value=None)`  
Convert Series to SparseSeries  
**Parameters** `kind` : {'block', 'integer'}  
`fill_value` : float, defaults to NaN (missing)  
**Returns** `sp` : SparseSeries

## 21.3 DataFrame

### 21.3.1 Attributes and underlying data

#### Axes

- **index**: row labels
- **columns**: column labels

<code>DataFrame.as_matrix([columns])</code>	Convert the frame to its Numpy-array matrix representation. Columns
<code>DataFrame.dtypes</code>	
<code>DataFrame.get_dtype_counts()</code>	
<code>DataFrame.values</code>	Convert the frame to its Numpy-array matrix representation. Columns
<code>DataFrame.axes</code>	
<code>DataFrame.ndim</code>	
<code>DataFrame.shape</code>	

## pandas.DataFrame.as\_matrix

`DataFrame.as_matrix(columns=None)`  
Convert the frame to its Numpy-array matrix representation. Columns are presented in sorted order unless a specific list of columns is provided.  
**Parameters** `columns` : array-like  
Specific column order  
**Returns** `values` : ndarray  
If the DataFrame is heterogeneous and contains booleans or objects, the result will be of dtype=object

## pandas.DataFrame.dtypes

`DataFrame.dtypes`

**pandas.DataFrame.get\_dtype\_counts**`DataFrame.get_dtype_counts()`**pandas.DataFrame.values**`DataFrame.values`

Convert the frame to its Numpy-array matrix representation. Columns are presented in sorted order unless a specific list of columns is provided.

**Parameters** `columns` : array-like

Specific column order

**Returns** `values` : ndarray

If the DataFrame is heterogeneous and contains booleans or objects, the result will be of dtype=object

**pandas.DataFrame.axes**`DataFrame.axes`**pandas.DataFrame.ndim**`DataFrame.ndim`**pandas.DataFrame.shape**`DataFrame.shape`**21.3.2 Conversion / Constructors**

<code>DataFrame.__init__([data, index, columns, ...])</code>	Two-dimensional size-mutable, potentially heterogeneous tabular data structure
<code>DataFrame.astype(dtype)</code>	Cast object to input numpy.dtype
<code>DataFrame.copy([deep])</code>	Make a copy of this object

**pandas.DataFrame.\_\_init\_\_**`DataFrame.__init__(data=None, index=None, columns=None, dtype=None, copy=False)`

Two-dimensional size-mutable, potentially heterogeneous tabular data structure with labeled axes (rows and columns). Arithmetic operations align on both row and column labels. Can be thought of as a dict-like container for Series objects. The primary pandas data structure

**Parameters** `data` : numpy ndarray (structured or homogeneous), dict, or DataFrame

Dict can contain Series, arrays, constants, or list-like objects

`index` : Index or array-like

Index to use for resulting frame. Will default to `np.arange(n)` if no indexing information part of input data and no index provided

**columns** : Index or array-like

Will default to `np.arange(n)` if not column labels provided

**dtype** : dtype, default None

Data type to force, otherwise infer

**copy** : boolean, default False

Copy data from inputs. Only affects DataFrame / 2d ndarray input

**See Also:**

**DataFrame.from\_records** constructor from tuples, also record arrays

**DataFrame.from\_dict** from dicts of Series, arrays, or dicts

**DataFrame.from\_csv** from CSV files

**DataFrame.from\_items** from sequence of (key, value) pairs

`read_csv`

### Examples

```
>>> d = {'col1': ts1, 'col2': ts2}
>>> df = DataFrame(data=d, index=index)
>>> df2 = DataFrame(np.random.randn(10, 5))
>>> df3 = DataFrame(np.random.randn(10, 5),
...                  columns=['a', 'b', 'c', 'd', 'e'])
```

## pandas.DataFrame.astype

**DataFrame.astype** (*dtype*)

Cast object to input `numpy.dtype`

**Parameters** **dtype** : `numpy.dtype` or Python type

**Returns** **casted** : type of caller

## pandas.DataFrame.copy

**DataFrame.copy** (*deep=True*)

Make a copy of this object

**Parameters** **deep** : boolean, default True

Make a deep copy, i.e. also copy data

**Returns** **copy** : type of caller

## 21.3.3 Indexing, iteration

---

`DataFrame.ix`

`DataFrame.insert`(loc, column, value) Insert column into DataFrame at specified location. Raises Exception if

Continued on next page

---

Table 21.23 – continued from previous page

<code>DataFrame.__iter__()</code>	Iterate over columns of the frame.
<code>DataFrame.iteritems()</code>	Iterator over (column, series) pairs
<code>DataFrame.pop(item)</code>	Return column and drop from frame.
<code>DataFrame.xs(key[, axis, level, copy])</code>	Returns a cross-section (row or column) from the DataFrame as a Series

**pandas.DataFrame.ix**`DataFrame.ix`**pandas.DataFrame.insert**`DataFrame.insert` (*loc, column, value*)

Insert column into DataFrame at specified location. Raises Exception if column is already contained in the DataFrame

**Parameters** `loc` : intMust have  $0 \leq \text{loc} \leq \text{len}(\text{columns})$ **column** : object**value** : int, Series, or array-like**pandas.DataFrame.\_\_iter\_\_**`DataFrame.__iter__()`

Iterate over columns of the frame.

**pandas.DataFrame.iteritems**`DataFrame.iteritems()`

Iterator over (column, series) pairs

**pandas.DataFrame.pop**`DataFrame.pop` (*item*)

Return column and drop from frame. Raise KeyError if not found.

**Returns** `column` : Series**pandas.DataFrame.xs**`DataFrame.xs` (*key, axis=0, level=None, copy=True*)

Returns a cross-section (row or column) from the DataFrame as a Series object. Defaults to returning a row (axis 0)

**Parameters** `key` : object

Some label contained in the index, or partially in a MultiIndex

**axis** : int, default 0

Axis to retrieve cross-section on

**copy** : boolean, default True

Whether to make a copy of the data

**Returns** **xs** : Series

### 21.3.4 Binary operator functions

<code>DataFrame.add(other[, axis, level, fill_value])</code>	Binary operator add with support to substitute a fill_value for missing data in
<code>DataFrame.div(other[, axis, level, fill_value])</code>	Binary operator divide with support to substitute a fill_value for missing data in
<code>DataFrame.mul(other[, axis, level, fill_value])</code>	Binary operator multiply with support to substitute a fill_value for missing data in
<code>DataFrame.sub(other[, axis, level, fill_value])</code>	Binary operator subtract with support to substitute a fill_value for missing data in
<code>DataFrame.radd(other[, axis, level, fill_value])</code>	Binary operator radd with support to substitute a fill_value for missing data in
<code>DataFrame.rdiv(other[, axis, level, fill_value])</code>	Binary operator rdivide with support to substitute a fill_value for missing data in
<code>DataFrame.rmul(other[, axis, level, fill_value])</code>	Binary operator rmultiply with support to substitute a fill_value for missing data in
<code>DataFrame.rsub(other[, axis, level, fill_value])</code>	Binary operator rsubtract with support to substitute a fill_value for missing data in
<code>DataFrame.combine(other, func[, fill_value])</code>	Add two DataFrame objects and do not propagate NaN values, so if for a
<code>DataFrame.combineAdd(other)</code>	Add two DataFrame objects and do not propagate
<code>DataFrame.combine_first(other)</code>	Combine two DataFrame objects and default to non-null values in frame
<code>DataFrame.combineMult(other)</code>	Multiply two DataFrame objects and do not propagate NaN values, so if

#### pandas.DataFrame.add

`DataFrame.add(other, axis='columns', level=None, fill_value=None)`

Binary operator add with support to substitute a fill\_value for missing data in one of the inputs

**Parameters** **other** : Series, DataFrame, or constant

**axis** : {0, 1, 'index', 'columns'}

For Series input, axis to match Series index on

**fill\_value** : None or float value, default None

Fill missing (NaN) values with this value. If both DataFrame locations are missing, the result will be missing

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**Returns** **result** : DataFrame

#### Notes

Mismatched indices will be unioned together

#### pandas.DataFrame.div

`DataFrame.div(other, axis='columns', level=None, fill_value=None)`

Binary operator divide with support to substitute a fill\_value for missing data in one of the inputs

**Parameters** **other** : Series, DataFrame, or constant

**axis** : {0, 1, 'index', 'columns'}

For Series input, axis to match Series index on

**fill\_value** : None or float value, default None

Fill missing (NaN) values with this value. If both DataFrame locations are missing, the result will be missing

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**Returns** **result** : DataFrame

#### Notes

Mismatched indices will be unioned together

### pandas.DataFrame.mul

DataFrame.**mul** (*other*, axis='columns', level=None, fill\_value=None)

Binary operator multiply with support to substitute a fill\_value for missing data in one of the inputs

**Parameters** **other** : Series, DataFrame, or constant

**axis** : {0, 1, 'index', 'columns'}

For Series input, axis to match Series index on

**fill\_value** : None or float value, default None

Fill missing (NaN) values with this value. If both DataFrame locations are missing, the result will be missing

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**Returns** **result** : DataFrame

#### Notes

Mismatched indices will be unioned together

### pandas.DataFrame.sub

DataFrame.**sub** (*other*, axis='columns', level=None, fill\_value=None)

Binary operator subtract with support to substitute a fill\_value for missing data in one of the inputs

**Parameters** **other** : Series, DataFrame, or constant

**axis** : {0, 1, 'index', 'columns'}

For Series input, axis to match Series index on

**fill\_value** : None or float value, default None

Fill missing (NaN) values with this value. If both DataFrame locations are missing, the result will be missing

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**Returns** **result** : DataFrame

#### Notes

Mismatched indices will be unioned together

### **pandas.DataFrame.radd**

`DataFrame.radd(other, axis='columns', level=None, fill_value=None)`

Binary operator radd with support to substitute a fill\_value for missing data in one of the inputs

**Parameters** **other** : Series, DataFrame, or constant

**axis** : {0, 1, 'index', 'columns'}

For Series input, axis to match Series index on

**fill\_value** : None or float value, default None

Fill missing (NaN) values with this value. If both DataFrame locations are missing, the result will be missing

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**Returns** **result** : DataFrame

#### Notes

Mismatched indices will be unioned together

### **pandas.DataFrame.rdiv**

`DataFrame.rdiv(other, axis='columns', level=None, fill_value=None)`

Binary operator rdivide with support to substitute a fill\_value for missing data in one of the inputs

**Parameters** **other** : Series, DataFrame, or constant

**axis** : {0, 1, 'index', 'columns'}

For Series input, axis to match Series index on

**fill\_value** : None or float value, default None

Fill missing (NaN) values with this value. If both DataFrame locations are missing, the result will be missing

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**Returns** **result** : DataFrame



### Notes

Mismatched indices will be unioned together

### pandas.DataFrame.rmul

`DataFrame.rmul` (*other*, *axis*='columns', *level*=None, *fill\_value*=None)

Binary operator rmultiply with support to substitute a *fill\_value* for missing data in one of the inputs

**Parameters** **other** : Series, DataFrame, or constant

**axis** : {0, 1, 'index', 'columns'}

For Series input, axis to match Series index on

**fill\_value** : None or float value, default None

Fill missing (NaN) values with this value. If both DataFrame locations are missing, the result will be missing

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**Returns** **result** : DataFrame

### Notes

Mismatched indices will be unioned together

### pandas.DataFrame.rsub

`DataFrame.rsub` (*other*, *axis*='columns', *level*=None, *fill\_value*=None)

Binary operator rsubtract with support to substitute a *fill\_value* for missing data in one of the inputs

**Parameters** **other** : Series, DataFrame, or constant

**axis** : {0, 1, 'index', 'columns'}

For Series input, axis to match Series index on

**fill\_value** : None or float value, default None

Fill missing (NaN) values with this value. If both DataFrame locations are missing, the result will be missing

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**Returns** **result** : DataFrame

### Notes

Mismatched indices will be unioned together

## pandas.DataFrame.combine

`DataFrame.combine` (*other*, *func*, *fill\_value=None*)

Add two DataFrame objects and do not propagate NaN values, so if for a (column, time) one frame is missing a value, it will default to the other frame's value (which might be NaN as well)

**Parameters** *other* : DataFrame

*func* : function

*fill\_value* : scalar value

**Returns** *result* : DataFrame

## pandas.DataFrame.combineAdd

`DataFrame.combineAdd` (*other*)

Add two DataFrame objects and do not propagate NaN values, so if for a (column, time) one frame is missing a value, it will default to the other frame's value (which might be NaN as well)

**Parameters** *other* : DataFrame

**Returns** DataFrame :

## pandas.DataFrame.combine\_first

`DataFrame.combine_first` (*other*)

Combine two DataFrame objects and default to non-null values in frame calling the method. Result index will be the union of the two indexes

**Parameters** *other* : DataFrame

**Returns** *combined* : DataFrame

### Examples

```
>>> a.combine_first(b)
a's values prioritized, use values from b to fill holes
```

## pandas.DataFrame.combineMult

`DataFrame.combineMult` (*other*)

Multiply two DataFrame objects and do not propagate NaN values, so if for a (column, time) one frame is missing a value, it will default to the other frame's value (which might be NaN as well)

**Parameters** *other* : DataFrame

**Returns** DataFrame :

## 21.3.5 Function application, GroupBy

---

<code>DataFrame.apply</code> ( <i>func</i> [, <i>axis</i> , <i>broadcast</i> , ...])	Applies function along input axis of DataFrame. Objects passed to
<code>DataFrame.applymap</code> ( <i>func</i> )	Apply a function to a DataFrame that is intended to operate

---

Continued on next page

Table 21.25 – continued from previous page

<code>DataFrame.groupby([by, axis, level, ...])</code>	Group series using mapper (dict or key function, apply given function)
--	--

**pandas.DataFrame.apply**

`DataFrame.apply(func, axis=0, broadcast=False, raw=False, args=(), **kws)`

Applies function along input axis of DataFrame. Objects passed to functions are Series objects having index either the DataFrame's index (axis=0) or the columns (axis=1). Return type depends on whether passed function aggregates

**Parameters** **func** : function

Function to apply to each column

**axis** : {0, 1}

0 : apply function to each column 1 : apply function to each row

**broadcast** : bool, default False

For aggregation functions, return object of same size with values propagated

**raw** : boolean, default False

If False, convert each row or column into a Series. If raw=True the passed function will receive ndarray objects instead. If you are just applying a NumPy reduction function this will achieve much better performance

**args** : tuple

Positional arguments to pass to function in addition to the array/series

**Additional keyword arguments will be passed as keywords to the function :**

**Returns** **applied** : Series or DataFrame

**Notes**

To apply a function elementwise, use `applymap`

**Examples**

```
>>> df.apply(numpy.sqrt) # returns DataFrame
>>> df.apply(numpy.sum, axis=0) # equiv to df.sum(0)
>>> df.apply(numpy.sum, axis=1) # equiv to df.sum(1)
```

**pandas.DataFrame.applymap**

`DataFrame.applymap(func)`

Apply a function to a DataFrame that is intended to operate elementwise, i.e. like doing `map(func, series)` for each series in the DataFrame

**Parameters** **func** : function

Python function, returns a single value from a single value

**Returns** **applied** : DataFrame

**pandas.DataFrame.groupby**

`DataFrame.groupby` (*by=None, axis=0, level=None, as\_index=True, sort=True, group\_keys=True*)

Group series using mapper (dict or key function, apply given function to group, return result as series) or by a series of columns

**Parameters** **by** : mapping function / list of functions, dict, Series, or tuple /

list of column names. Called on each element of the object index to determine the groups. If a dict or Series is passed, the Series or dict VALUES will be used to determine the groups

**axis** : int, default 0

**level** : int, level name, or sequence of such, default None

If the axis is a MultiIndex (hierarchical), group by a particular level or levels

**as\_index** : boolean, default True

For aggregated output, return object with group labels as the index. Only relevant for DataFrame input. `as_index=False` is effectively “SQL-style” grouped output

**sort** : boolean, default True

Sort group keys. Get better performance by turning this off

**group\_keys** : boolean, default True

When calling `apply`, add group keys to index to identify pieces

**Returns** **GroupBy object** :

**Examples**

```
# DataFrame result >>> data.groupby(func, axis=0).mean()
# DataFrame result >>> data.groupby(['col1', 'col2'])['col3'].mean()
# DataFrame with hierarchical index >>> data.groupby(['col1', 'col2']).mean()
```

**21.3.6 Computations / Descriptive Stats**

<code>DataFrame.clip([upper, lower])</code>	Trim values at input threshold(s)
<code>DataFrame.clip_lower(threshold)</code>	Trim values below threshold
<code>DataFrame.clip_upper(threshold)</code>	Trim values above threshold
<code>DataFrame.corr([method])</code>	Compute pairwise correlation of columns, excluding NA/null values
<code>DataFrame.corrwith(other[, axis, drop])</code>	Compute pairwise correlation between rows or columns of two DataFrame
<code>DataFrame.count([axis, level, numeric_only])</code>	Return Series with number of non-NA/null observations over requested
<code>DataFrame.cumprod([axis, skipna])</code>	Return cumulative product over requested axis as DataFrame
<code>DataFrame.cumsum([axis, skipna])</code>	Return DataFrame of cumulative sums over requested axis.
<code>DataFrame.describe([percentile_width])</code>	Generate various summary statistics of each column, excluding
<code>DataFrame.diff([periods])</code>	1st discrete difference of object
<code>DataFrame.mad([axis, skipna, level])</code>	Return mean absolute deviation over requested axis.
<code>DataFrame.max([axis, skipna, level])</code>	Return maximum over requested axis.
<code>DataFrame.mean([axis, skipna, level])</code>	Return mean over requested axis.
<code>DataFrame.median([axis, skipna, level])</code>	Return median over requested axis.

Continued on next page

Table 21.26 – continued from previous page

<code>DataFrame.min([axis, skipna, level])</code>	Return minimum over requested axis.
<code>DataFrame.prod([axis, skipna, level])</code>	Return product over requested axis.
<code>DataFrame.quantile([q, axis])</code>	Return values at the given quantile over requested axis, a la
<code>DataFrame.skew([axis, skipna, level])</code>	Return unbiased skewness over requested axis.
<code>DataFrame.sum([axis, numeric_only, skipna, ...])</code>	Return sum over requested axis.
<code>DataFrame.std([axis, skipna, level, ddof])</code>	Return standard deviation over requested axis.
<code>DataFrame.var([axis, skipna, level, ddof])</code>	Return variance over requested axis.

**pandas.DataFrame.clip**

`DataFrame.clip` (*upper=None, lower=None*)

Trim values at input threshold(s)

**Parameters** `lower` : float, default None

`upper` : float, default None

**Returns** `clipped` : DataFrame

**pandas.DataFrame.clip\_lower**

`DataFrame.clip_lower` (*threshold*)

Trim values below threshold

**Returns** `clipped` : DataFrame

**pandas.DataFrame.clip\_upper**

`DataFrame.clip_upper` (*threshold*)

Trim values above threshold

**Returns** `clipped` : DataFrame

**pandas.DataFrame.corr**

`DataFrame.corr` (*method='pearson'*)

Compute pairwise correlation of columns, excluding NA/null values

**Parameters** `method` : {'pearson', 'kendall', 'spearman'}

pearson : standard correlation coefficient kendall : Kendall Tau correlation coefficient

spearman : Spearman rank correlation

**Returns** `y` : DataFrame

**pandas.DataFrame.corrwith**

`DataFrame.corrwith` (*other, axis=0, drop=False*)

Compute pairwise correlation between rows or columns of two DataFrame objects.

**Parameters** `other` : DataFrame

`axis` : {0, 1}

0 to compute column-wise, 1 for row-wise

**drop** : boolean, default False

Drop missing indices from result, default returns union of all

**Returns** **correls** : Series

## **pandas.DataFrame.count**

`DataFrame.count` (*axis=0, level=None, numeric\_only=False*)

Return Series with number of non-NA/null observations over requested axis. Works with non-floating point data as well (detects NaN and None)

**Parameters** **axis** : {0, 1}

0 for row-wise, 1 for column-wise

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a DataFrame

**numeric\_only** : boolean, default False

Include only float, int, boolean data

**Returns** **count** : Series (or DataFrame if level specified)

## **pandas.DataFrame.cumprod**

`DataFrame.cumprod` (*axis=None, skipna=True*)

Return cumulative product over requested axis as DataFrame

**Parameters** **axis** : {0, 1}

0 for row-wise, 1 for column-wise

**skipna** : boolean, default True

Exclude NA/null values. If an entire row/column is NA, the result will be NA

**Returns** **y** : DataFrame

## **pandas.DataFrame.cumsum**

`DataFrame.cumsum` (*axis=None, skipna=True*)

Return DataFrame of cumulative sums over requested axis.

**Parameters** **axis** : {0, 1}

0 for row-wise, 1 for column-wise

**skipna** : boolean, default True

Exclude NA/null values. If an entire row/column is NA, the result will be NA

**Returns** **y** : DataFrame

### pandas.DataFrame.describe

DataFrame.**describe** (*percentile\_width=50*)

Generate various summary statistics of each column, excluding NaN values. These include: count, mean, std, min, max, and lower%/50%/upper% percentiles

**Parameters** **percentile\_width** : float, optional

width of the desired uncertainty interval, default is 50, which corresponds to lower=25, upper=75

**Returns** **DataFrame of summary statistics** :

### pandas.DataFrame.diff

DataFrame.**diff** (*periods=1*)

1st discrete difference of object

**Parameters** **periods** : int, default 1

Periods to shift for forming difference

**Returns** **difff** : DataFrame

### pandas.DataFrame.mad

DataFrame.**mad** (*axis=0, skipna=True, level=None*)

Return mean absolute deviation over requested axis. NA/null values are excluded

**Parameters** **axis** : {0, 1}

0 for row-wise, 1 for column-wise

**skipna** : boolean, default True

Exclude NA/null values. If an entire row/column is NA, the result will be NA

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a DataFrame

**Returns** **mad** : Series (or DataFrame if level specified)

### pandas.DataFrame.max

DataFrame.**max** (*axis=0, skipna=True, level=None*)

Return maximum over requested axis. NA/null values are excluded

**Parameters** **axis** : {0, 1}

0 for row-wise, 1 for column-wise

**skipna** : boolean, default True

Exclude NA/null values. If an entire row/column is NA, the result will be NA

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a DataFrame

**Returns** **max** : Series (or DataFrame if level specified)

### **pandas.DataFrame.mean**

`DataFrame.mean` (*axis=0, skipna=True, level=None*)

Return mean over requested axis. NA/null values are excluded

**Parameters** **axis** : {0, 1}

0 for row-wise, 1 for column-wise

**skipna** : boolean, default True

Exclude NA/null values. If an entire row/column is NA, the result will be NA

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a DataFrame

**Returns** **mean** : Series (or DataFrame if level specified)

### **pandas.DataFrame.median**

`DataFrame.median` (*axis=0, skipna=True, level=None*)

Return median over requested axis. NA/null values are excluded

**Parameters** **axis** : {0, 1}

0 for row-wise, 1 for column-wise

**skipna** : boolean, default True

Exclude NA/null values. If an entire row/column is NA, the result will be NA

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a DataFrame

**Returns** **median** : Series (or DataFrame if level specified)

### **pandas.DataFrame.min**

`DataFrame.min` (*axis=0, skipna=True, level=None*)

Return minimum over requested axis. NA/null values are excluded

**Parameters** **axis** : {0, 1}

0 for row-wise, 1 for column-wise

**skipna** : boolean, default True

Exclude NA/null values. If an entire row/column is NA, the result will be NA

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a DataFrame

**Returns** **min** : Series (or DataFrame if level specified)



**pandas.DataFrame.prod**`DataFrame.prod(axis=0, skipna=True, level=None)`

Return product over requested axis. NA/null values are treated as 1

**Parameters** `axis` : {0, 1}

0 for row-wise, 1 for column-wise

**skipna** : boolean, default True

Exclude NA/null values. If an entire row/column is NA, the result will be NA

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a DataFrame

**Returns** `product` : Series (or DataFrame if level specified)**pandas.DataFrame.quantile**`DataFrame.quantile(q=0.5, axis=0)`Return values at the given quantile over requested axis, a la `scoreatpercentile` in `scipy.stats`**Parameters** `q` : quantile, default 0.5 (50% quantile)

0 &lt;= q &lt;= 1

**axis** : {0, 1}

0 for row-wise, 1 for column-wise

**Returns** `quantiles` : Series**pandas.DataFrame.skew**`DataFrame.skew(axis=0, skipna=True, level=None)`

Return unbiased skewness over requested axis. NA/null values are excluded

**Parameters** `axis` : {0, 1}

0 for row-wise, 1 for column-wise

**skipna** : boolean, default True

Exclude NA/null values. If an entire row/column is NA, the result will be NA

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a DataFrame

**Returns** `skew` : Series (or DataFrame if level specified)**pandas.DataFrame.sum**`DataFrame.sum(axis=0, numeric_only=None, skipna=True, level=None)`

Return sum over requested axis. NA/null values are excluded

**Parameters** `axis` : {0, 1}

0 for row-wise, 1 for column-wise

**skipna** : boolean, default True

Exclude NA/null values. If an entire row/column is NA, the result will be NA

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a DataFrame

**numeric\_only** : boolean, default None

Include only float, int, boolean data. If None, will attempt to use everything, then use only numeric data

**Returns** **sum** : Series (or DataFrame if level specified)

### **pandas.DataFrame.std**

`DataFrame.std(axis=0, skipna=True, level=None, ddof=1)`

Return standard deviation over requested axis. NA/null values are excluded

**Parameters** **axis** : {0, 1}

0 for row-wise, 1 for column-wise

**skipna** : boolean, default True

Exclude NA/null values. If an entire row/column is NA, the result will be NA

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a DataFrame

**Returns** **std** : Series (or DataFrame if level specified)

### **pandas.DataFrame.var**

`DataFrame.var(axis=0, skipna=True, level=None, ddof=1)`

Return variance over requested axis. NA/null values are excluded

**Parameters** **axis** : {0, 1}

0 for row-wise, 1 for column-wise

**skipna** : boolean, default True

Exclude NA/null values. If an entire row/column is NA, the result will be NA

**level** : int, default None

If the axis is a MultiIndex (hierarchical), count along a particular level, collapsing into a DataFrame

**Returns** **var** : Series (or DataFrame if level specified)

---

Continued on next page

Table 21.27 – continued from previous page

### 21.3.7 Reindexing / Selection / Label manipulation

<code>DataFrame.add_prefix(prefix)</code>	Concatenate prefix string with panel items names.
<code>DataFrame.add_suffix(suffix)</code>	Concatenate suffix string with panel items names
<code>DataFrame.align(other[, join, axis, level, ...])</code>	Align two DataFrame object on their index and columns with the
<code>DataFrame.drop(labels[, axis, level])</code>	Return new object with labels in requested axis removed
<code>DataFrame.filter([items, like, regex])</code>	Restrict frame's columns to set of items or wildcard
<code>DataFrame.reindex([index, columns, method, ...])</code>	Conform DataFrame to new index with optional filling logic, placing
<code>DataFrame.reindex_like(other[, method, ...])</code>	Reindex DataFrame to match indices of another DataFrame, optionally
<code>DataFrame.rename([index, columns, copy, inplace])</code>	Alter index and / or columns using input function or functions.
<code>DataFrame.select(crit[, axis])</code>	Return data corresponding to axis labels matching criteria
<code>DataFrame.take(indices[, axis])</code>	Analogous to ndarray.take, return DataFrame corresponding to requested
<code>DataFrame.truncate([before, after, copy])</code>	Function truncate a sorted DataFrame / Series before and/or after
<code>DataFrame.head([n])</code>	Returns first n rows of DataFrame
<code>DataFrame.tail([n])</code>	Returns last n rows of DataFrame

#### pandas.DataFrame.add\_prefix

`DataFrame.add_prefix` (*prefix*)

Concatenate prefix string with panel items names.

**Parameters** `prefix` : string

**Returns** `with_prefix` : type of caller

#### pandas.DataFrame.add\_suffix

`DataFrame.add_suffix` (*suffix*)

Concatenate suffix string with panel items names

**Parameters** `suffix` : string

**Returns** `with_suffix` : type of caller

#### pandas.DataFrame.align

`DataFrame.align` (*other*, *join*='outer', *axis*=None, *level*=None, *copy*=True, *fill\_value*=nan, *method*=None, *limit*=None, *fill\_axis*=0)

Align two DataFrame object on their index and columns with the specified join method for each axis Index

**Parameters** `other` : DataFrame or Series

**join** : { 'outer', 'inner', 'left', 'right' }, default 'outer'

**axis** : {0, 1, None}, default None

Align on index (0), columns (1), or both (None)

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**copy** : boolean, default True

Always returns new objects. If `copy=False` and no reindexing is required then original objects are returned.

**fill\_value** : scalar, default `np.NaN`

Value to use for missing values. Defaults to `NaN`, but can be any “compatible” value

**method** : str, default `None`

**limit** : int, default `None`

**fill\_axis** : {0, 1}, default 0

Filling axis, method and limit

**Returns** (**left, right**) : (DataFrame, type of other)

Aligned objects

### **pandas.DataFrame.drop**

`DataFrame.drop(labels, axis=0, level=None)`

Return new object with labels in requested axis removed

**Parameters** **labels** : array-like

**axis** : int

**level** : int or name, default `None`

For MultiIndex

**Returns** **dropped** : type of caller

### **pandas.DataFrame.filter**

`DataFrame.filter(items=None, like=None, regex=None)`

Restrict frame’s columns to set of items or wildcard

**Parameters** **items** : list-like

List of columns to restrict to (must not all be present)

**like** : string

Keep columns where “arg in col == True”

**regex** : string (regular expression)

Keep columns with `re.search(regex, col) == True`

**Returns** **DataFrame with filtered columns** :

#### **Notes**

Arguments are mutually exclusive, but this is not checked for

**pandas.DataFrame.reindex**

`DataFrame.reindex(index=None, columns=None, method=None, level=None, fill_value=nan, limit=None, copy=True)`

Conform DataFrame to new index with optional filling logic, placing NA/NaN in locations having no value in the previous index. A new object is produced unless the new index is equivalent to the current one and `copy=False`

**Parameters** **index** : array-like, optional

New labels / index to conform to. Preferably an Index object to avoid duplicating data

**columns** : array-like, optional

Same usage as index argument

**method** : { 'backfill', 'bfill', 'pad', 'ffill', None }, default None

Method to use for filling holes in reindexed DataFrame `pad` / `ffill`: propagate last valid observation forward to next valid `backfill` / `bfill`: use NEXT valid observation to fill gap

**copy** : boolean, default True

Return a new object, even if the passed indexes are the same

**level** : int or name

Broadcast across a level, matching Index values on the passed MultiIndex level

**fill\_value** : scalar, default np.NaN

Value to use for missing values. Defaults to NaN, but can be any “compatible” value

**limit** : int, default None

Maximum size gap to forward or backward fill

**Returns** **reindexed** : same type as calling instance

**Examples**

```
>>> df.reindex(index=[date1, date2, date3], columns=['A', 'B', 'C'])
```

**pandas.DataFrame.reindex\_like**

`DataFrame.reindex_like(other, method=None, copy=True, limit=None)`

Reindex DataFrame to match indices of another DataFrame, optionally with filling logic

**Parameters** **other** : DataFrame

**method** : string or None

**copy** : boolean, default True

**limit** : int, default None

Maximum size gap to forward or backward fill

**Returns** **reindexed** : DataFrame

## Notes

Like calling `s.reindex(index=other.index, columns=other.columns, method=...)`

## pandas.DataFrame.rename

`DataFrame.rename` (*index=None, columns=None, copy=True, inplace=False*)

Alter index and / or columns using input function or functions. Function / dict values must be unique (1-to-1). Labels not contained in a dict / Series will be left as-is.

**Parameters** **index** : dict-like or function, optional

Transformation to apply to index values

**columns** : dict-like or function, optional

Transformation to apply to column values

**copy** : boolean, default True

Also copy underlying data

**inplace** : boolean, default False

Whether to return a new DataFrame. If True then value of copy is ignored.

**Returns** **renamed** : DataFrame (new object)

**See Also:**

[`Series.rename`](#)

## pandas.DataFrame.select

`DataFrame.select` (*crit, axis=0*)

Return data corresponding to axis labels matching criteria

**Parameters** **crit** : function

To be called on each index (label). Should return True or False

**axis** : int

**Returns** **selection** : type of caller

## pandas.DataFrame.take

`DataFrame.take` (*indices, axis=0*)

Analogous to `ndarray.take`, return DataFrame corresponding to requested indices along an axis

**Parameters** **indices** : list / array of ints

**axis** : {0, 1}

**Returns** **taken** : DataFrame

**pandas.DataFrame.truncate**

`DataFrame.truncate` (*before=None, after=None, copy=True*)

Function truncate a sorted DataFrame / Series before and/or after some particular dates.

**Parameters** **before** : date

Truncate before date

**after** : date

Truncate after date

**Returns** **truncated** : type of caller

**pandas.DataFrame.head**

`DataFrame.head` (*n=5*)

Returns first n rows of DataFrame

**pandas.DataFrame.tail**

`DataFrame.tail` (*n=5*)

Returns last n rows of DataFrame

**21.3.8 Missing data handling**

<code>DataFrame.dropna</code> ([axis, how, thresh, subset])	Return object with labels on given axis omitted where alternately any
<code>DataFrame.fillna</code> ([value, method, axis, ...])	Fill NA/NaN values using the specified method

**pandas.DataFrame.dropna**

`DataFrame.dropna` (*axis=0, how='any', thresh=None, subset=None*)

Return object with labels on given axis omitted where alternately any or all of the data are missing

**Parameters** **axis** : {0, 1}

**how** : {'any', 'all'}

any : if any NA values are present, drop that label  
all : if all values are NA, drop that label

**thresh** : int, default None

int value : require that many non-NA values

**subset** : array-like

Labels along other axis to consider, e.g. if you are dropping rows these would be a list of columns to include

**Returns** **dropped** : DataFrame

**pandas.DataFrame.fillna**

`DataFrame.fillna` (*value=None, method='pad', axis=0, inplace=False, limit=None*)

Fill NA/NaN values using the specified method

**Parameters** **method** : { 'backfill', 'bfill', 'pad', 'ffill', None }, default 'pad'

Method to use for filling holes in reindexed Series `pad` / `ffill`: propagate last valid observation forward to next valid `backfill` / `bfill`: use NEXT valid observation to fill gap

**value** : scalar or dict

Value to use to fill holes (e.g. 0), alternately a dict of values specifying which value to use for each column (columns not in the dict will not be filled)

**axis** : {0, 1}, default 0

0: fill column-by-column 1: fill row-by-row

**inplace** : boolean, default False

If True, fill the DataFrame in place. Note: this will modify any other views on this DataFrame, like if you took a no-copy slice of an existing DataFrame, for example a column in a DataFrame. Returns a reference to the filled object, which is self if `inplace=True`

**limit** : int, default None

Maximum size gap to forward or backward fill

**Returns** **filled** : DataFrame

**See Also:**

`reindex`, `asfreq`

### 21.3.9 Reshaping, sorting, transposing

<code>DataFrame.sort_index([axis, by, ascending, ...])</code>	Sort DataFrame either by labels (along either axis) or by the values in
<code>DataFrame.delevel(*args, **kwargs)</code>	
<code>DataFrame.pivot([index, columns, values])</code>	Reshape data (produce a “pivot” table) based on column values.
<code>DataFrame.sortlevel([level, axis, ascending])</code>	Sort multilevel index by chosen axis and primary level.
<code>DataFrame.swaplevel(i, j[, axis])</code>	Swap levels i and j in a MultiIndex on a particular axis
<code>DataFrame.stack([level, dropna])</code>	Pivot a level of the (possibly hierarchical) column labels, returning a
<code>DataFrame.unstack([level])</code>	Pivot a level of the (necessarily hierarchical) index labels, returning
<code>DataFrame.T</code>	Returns a DataFrame with the rows/columns switched. If the DataFrame is
<code>DataFrame.transpose()</code>	Returns a DataFrame with the rows/columns switched. If the DataFrame is

**pandas.DataFrame.sort\_index**

`DataFrame.sort_index` (*axis=0, by=None, ascending=True, inplace=False*)

Sort DataFrame either by labels (along either axis) or by the values in a column

**Parameters** **axis** : {0, 1}

Sort index/rows versus columns

**by** : object

Column name(s) in frame. Accepts a column name or a list or tuple for a nested sort.



**ascending** : boolean, default True

Sort ascending vs. descending

**inplace** : boolean, default False

Sort the DataFrame without creating a new instance

**Returns** **sorted** : DataFrame

### **pandas.DataFrame.delevel**

DataFrame.**delevel** (\*args, \*\*kwargs)

### **pandas.DataFrame.pivot**

DataFrame.**pivot** (index=None, columns=None, values=None)

Reshape data (produce a “pivot” table) based on column values. Uses unique values from index / columns to form axes and return either DataFrame or Panel, depending on whether you request a single value column (DataFrame) or all columns (Panel)

**Parameters** **index** : string or object

Column name to use to make new frame’s index

**columns** : string or object

Column name to use to make new frame’s columns

**values** : string or object, optional

Column name to use for populating new frame’s values

**Returns** **pivoted** : DataFrame

If no values column specified, will have hierarchically indexed columns

### **Notes**

For finer-tuned control, see hierarchical indexing documentation along with the related stack/unstack methods

### **Examples**

```
>>> df
   foo  bar  baz
0  one   A   1.
1  one   B   2.
2  one   C   3.
3  two   A   4.
4  two   B   5.
5  two   C   6.

>>> df.pivot('foo', 'bar', 'baz')
   A  B  C
one 1  2  3
two 4  5  6
```

```
>>> df.pivot('foo', 'bar')['baz']
      A  B  C
one  1  2  3
two  4  5  6
```

## pandas.DataFrame.sortlevel

DataFrame.**sortlevel** (*level=0, axis=0, ascending=True*)

Sort multilevel index by chosen axis and primary level. Data will be lexicographically sorted by the chosen level followed by the other levels (in order)

**Parameters** **level** : int

**axis** : {0, 1}

**ascending** : bool, default True

**Returns** **sorted** : DataFrame

## pandas.DataFrame.swaplevel

DataFrame.**swaplevel** (*i, j, axis=0*)

Swap levels i and j in a MultiIndex on a particular axis

**Returns** **swapped** : type of caller (new object)

## pandas.DataFrame.stack

DataFrame.**stack** (*level=-1, dropna=True*)

Pivot a level of the (possibly hierarchical) column labels, returning a DataFrame (or Series in the case of an object with a single level of column labels) having a hierarchical index with a new inner-most level of row labels.

**Parameters** **level** : int, string, or list of these, default last level

Level(s) to stack, can pass level name

**dropna** : boolean, default True

Whether to drop rows in the resulting Frame/Series with no valid values

**Returns** **stacked** : DataFrame or Series

## Examples

```
>>> s
      a  b
one  1.  2.
two  3.  4.

>>> s.stack()
one a    1
   b    2
two a    3
   b    4
```

**pandas.DataFrame.unstack**`DataFrame.unstack (level=-1)`

Pivot a level of the (necessarily hierarchical) index labels, returning a DataFrame having a new level of column labels whose inner-most level consists of the pivoted index labels. If the index is not a MultiIndex, the output will be a Series (the analogue of stack when the columns are not a MultiIndex)

**Parameters** **level** : int, string, or list of these, default last level

Level(s) of index to unstack, can pass level name

**Returns** **unstacked** : DataFrame or Series

**Examples**

```
>>> s
one  a    1.
one  b    2.
two  a    3.
two  b    4.

>>> s.unstack(level=-1)
      a    b
one  1.  2.
two  3.  4.

>>> df = s.unstack(level=0)
>>> df
      one  two
a    1.   2.
b    3.   4.

>>> df.unstack()
one  a    1.
     b    3.
two  a    2.
     b    4.
```

**pandas.DataFrame.T**`DataFrame.T`

Returns a DataFrame with the rows/columns switched. If the DataFrame is homogeneously-typed, the data is not copied

**pandas.DataFrame.transpose**`DataFrame.transpose()`

Returns a DataFrame with the rows/columns switched. If the DataFrame is homogeneously-typed, the data is not copied

**21.3.10 Combining / joining / merging**

<code>DataFrame.join(other[, on, how, lsuffix, ...])</code>	Join columns with other DataFrame either on index or on a key
<code>DataFrame.merge(right[, how, on, left_on, ...])</code>	Merge DataFrame objects by performing a database-style join operation by
<code>DataFrame.append(other[, ignore_index, ...])</code>	Append columns of other to end of this frame's columns and index, returning a

## pandas.DataFrame.join

`DataFrame.join(other, on=None, how='left', lsuffix='', rsuffix='', sort=False)`

Join columns with other DataFrame either on index or on a key column. Efficiently Join multiple DataFrame objects by index at once by passing a list.

**Parameters** **other** : DataFrame, Series with name field set, or list of DataFrame

Index should be similar to one of the columns in this one. If a Series is passed, its name attribute must be set, and that will be used as the column name in the resulting joined DataFrame

**on** : column name, tuple/list of column names, or array-like

Column(s) to use for joining, otherwise join on index. If multiples columns given, the passed DataFrame must have a MultiIndex. Can pass an array as the join key if not already contained in the calling DataFrame. Like an Excel VLOOKUP operation

**how** : { 'left', 'right', 'outer', 'inner' }

How to handle indexes of the two objects. Default: 'left' for joining on index, None otherwise \* left: use calling frame's index \* right: use input frame's index \* outer: form union of indexes \* inner: use intersection of indexes

**lsuffix** : string

Suffix to use from left frame's overlapping columns

**rsuffix** : string

Suffix to use from right frame's overlapping columns

**sort** : boolean, default False

Order result DataFrame lexicographically by the join key. If False, preserves the index order of the calling (left) DataFrame

**Returns** **joined** : DataFrame

### Notes

on, lsuffix, and rsuffix options are not supported when passing a list of DataFrame objects

## pandas.DataFrame.merge

`DataFrame.merge(right, how='inner', on=None, left_on=None, right_on=None, left_index=False, right_index=False, sort=True, suffixes=('_x', '_y'), copy=True)`

Merge DataFrame objects by performing a database-style join operation by columns or indexes.

If joining columns on columns, the DataFrame indexes *will be ignored*. Otherwise if joining indexes on indexes or indexes on a column or columns, the index will be passed on.

**Parameters** **right** : DataFrame

**how** : { 'left', 'right', 'outer', 'inner' }, default 'inner'

- left: use only keys from left frame (SQL: left outer join)
- right: use only keys from right frame (SQL: right outer join)
- outer: use union of keys from both frames (SQL: full outer join)
- inner: use intersection of keys from both frames (SQL: inner join)

**on** : label or list

Field names to join on. Must be found in both DataFrames.

**left\_on** : label or list, or array-like

Field names to join on in left DataFrame. Can be a vector or list of vectors of the length of the DataFrame to use a particular vector as the join key instead of columns

**right\_on** : label or list, or array-like

Field names to join on in right DataFrame or vector/list of vectors per left\_on docs

**left\_index** : boolean, default True

Use the index from the left DataFrame as the join key(s). If it is a MultiIndex, the number of keys in the other DataFrame (either the index or a number of columns) must match the number of levels

**right\_index** : boolean, default True

Use the index from the right DataFrame as the join key. Same caveats as left\_index

**sort** : boolean, default True

Sort the join keys lexicographically in the result DataFrame

**suffixes** : 2-length sequence (tuple, list, ...)

Suffix to apply to overlapping column names in the left and right side, respectively

**copy** : boolean, default True

If False, do not copy data unnecessarily

**Returns** **merged** : DataFrame

### Examples

```
>>> A
   lkey value
0  foo   1
1  bar   2
2  baz   3
3  foo   4

>>> B
   rkey value
0  foo   5
1  bar   6
2  qux   7
3  bar   8

>>> merge(A, B, left_on='lkey', right_on='rkey', how='outer')
   lkey  value_x  rkey  value_y
0  bar     2     bar     6
1  bar     2     bar     8
2  baz     3    NaN    NaN
3  foo     1     foo     5
4  foo     4     foo     5
5  NaN    NaN    qux     7
```

## pandas.DataFrame.append

`DataFrame.append(other, ignore_index=False, verify_integrity=False)`

Append columns of other to end of this frame's columns and index, returning a new object. Columns not in this frame are added as new columns.

**Parameters** **other** : DataFrame or list of Series/dict-like objects

**ignore\_index** : boolean, default False

If True do not use the index labels. Useful for gluing together record arrays

**verify\_integrity** : boolean, default False

If True, raise Exception on creating index with duplicates

**Returns** **appended** : DataFrame

### Notes

If a list of dict is passed and the keys are all contained in the DataFrame's index, the order of the columns in the resulting DataFrame will be unchanged

## 21.3.11 Time series-related

<code>DataFrame.asfreq(freq[, method, how])</code>	Convert all TimeSeries inside to specified frequency using DateOffset
<code>DataFrame.shift([periods, freq])</code>	Shift the index of the DataFrame by desired number of periods with an
<code>DataFrame.first_valid_index()</code>	Return label for first non-NA/null value
<code>DataFrame.last_valid_index()</code>	Return label for last non-NA/null value

## pandas.DataFrame.asfreq

`DataFrame.asfreq(freq, method=None, how=None)`

Convert all TimeSeries inside to specified frequency using DateOffset objects. Optionally provide fill method to pad/backfill missing values.

**Parameters** **freq** : DateOffset object, or string

**method** : { 'backfill', 'bfill', 'pad', 'ffill', None }

Method to use for filling holes in reindexed Series pad / ffill: propagate last valid observation forward to next valid backfill / bfill: use NEXT valid observation to fill methdo

**how** : { 'start', 'end' }, default end

For PeriodIndex only, see PeriodIndex.asfreq

**Returns** **converted** : type of caller

## pandas.DataFrame.shift

`DataFrame.shift(periods=1, freq=None, **kws)`

Shift the index of the DataFrame by desired number of periods with an optional time freq

**Parameters** **periods** : int

Number of periods to move, can be positive or negative

**freq** : DateOffset, timedelta, or time rule string, optional

Increment to use from datetools module or time rule (e.g. 'EOM')

**Returns** **shifted** : DataFrame

#### Notes

If freq is specified then the index values are shifted but the data is not realigned

#### pandas.DataFrame.first\_valid\_index

DataFrame.**first\_valid\_index**()

Return label for first non-NA/null value

#### pandas.DataFrame.last\_valid\_index

DataFrame.**last\_valid\_index**()

Return label for last non-NA/null value

### 21.3.12 Plotting

---

DataFrame.hist(data[, grid, xlabelsize, ...])	Draw Histogram the DataFrame's series using matplotlib / pylab.
---	---

---

DataFrame.plot([frame, x, y, subplots, ...])	Make line or bar plot of DataFrame's series with the index on the x-axis
--	--

---

#### pandas.DataFrame.hist

DataFrame.**hist** (data, grid=True, xlabelsize=None, xrot=None, ylabelsize=None, yrot=None, ax=None, sharex=False, sharey=False, \*\*kws)

Draw Histogram the DataFrame's series using matplotlib / pylab.

**Parameters** **grid** : boolean, default True

Whether to show axis grid lines

**xlabelsize** : int, default None

If specified changes the x-axis label size

**xrot** : float, default None

rotation of x axis labels

**ylabelsize** : int, default None

If specified changes the y-axis label size

**yrot** : float, default None

rotation of y axis labels

**ax** : matplotlib axes object, default None

**sharex** : bool, if True, the X axis will be shared amongst all subplots.

**sharey** : bool, if True, the Y axis will be shared amongst all subplots.

**kwds** : other plotting keyword arguments

To be passed to hist function

## pandas.DataFrame.plot

`DataFrame.plot` (*frame=None, x=None, y=None, subplots=False, sharex=True, sharey=False, use\_index=True, figsize=None, grid=False, legend=True, rot=None, ax=None, style=None, title=None, xlim=None, ylim=None, logy=False, xticks=None, yticks=None, kind='line', sort\_columns=False, fontsize=None, secondary\_y=False, \*\*kwds*)

Make line or bar plot of DataFrame's series with the index on the x-axis using matplotlib / pylab.

**Parameters** **x** : int or str, default None

**y** : int or str, default None

Allows plotting of one column versus another

**subplots** : boolean, default False

Make separate subplots for each time series

**sharex** : boolean, default True

In case subplots=True, share x axis

**sharey** : boolean, default False

In case subplots=True, share y axis

**use\_index** : boolean, default True

Use index as ticks for x axis

**stacked** : boolean, default False

If True, create stacked bar plot. Only valid for DataFrame input

**sort\_columns**: boolean, default False :

Sort column names to determine plot ordering

**title** : string

Title to use for the plot

**grid** : boolean, default True

Axis grid lines

**legend** : boolean, default True

Place legend on axis subplots

**ax** : matplotlib axis object, default None

**style** : list or dict

matplotlib line style per column

**kind** : { 'line', 'bar', 'barh' }

bar : vertical bar plot barh : horizontal bar plot

**logy** : boolean, default False

For line plots, use log scaling on y axis



**xticks** : sequence

Values to use for the xticks

**yticks** : sequence

Values to use for the yticks

**xlim** : 2-tuple/list

**ylim** : 2-tuple/list

**rot** : int, default None

Rotation for ticks

**secondary\_y** : boolean or sequence, default False

Whether to plot on the secondary y-axis. If dict then can select which columns to plot on secondary y-axis

**kwds** : keywords

Options to pass to matplotlib plotting method

**Returns** **ax\_or\_axes** : matplotlib.AxesSubplot or list of them

### 21.3.13 Serialization / IO / Conversion

<code>DataFrame.from_csv(path[, header, sep, ...])</code>	Read delimited file into DataFrame
<code>DataFrame.from_records(data[, index, ...])</code>	Convert structured or record ndarray to DataFrame
<code>DataFrame.to_csv(path_or_buf[, sep, na_rep, ...])</code>	Write DataFrame to a comma-separated values (csv) file
<code>DataFrame.to_excel(excel_writer[, ...])</code>	Write DataFrame to a excel sheet
<code>DataFrame.to_dict([outtype])</code>	Convert DataFrame to dictionary.
<code>DataFrame.to_records([index])</code>	Convert DataFrame to record array. Index will be put in the
<code>DataFrame.to_sparse([fill_value, kind])</code>	Convert to SparseDataFrame
<code>DataFrame.to_string([buf, columns, ...])</code>	Render a DataFrame to a console-friendly tabular output.
<code>DataFrame.save(path)</code>	
<code>DataFrame.load(path)</code>	
<code>DataFrame.info([verbose, buf])</code>	Concise summary of a DataFrame, used in <code>__repr__</code> when very large.

#### pandas.DataFrame.from\_csv

**classmethod** `DataFrame.from_csv` (*path*, *header=0*, *sep=''*, *index\_col=0*, *parse\_dates=True*, *encoding=None*)

Read delimited file into DataFrame

**Parameters** **path** : string file path or file handle / StringIO

**header** : int, default 0

Row to use at header (skip prior rows)

**sep** : string, default ','

Field delimiter

**index\_col** : int or sequence, default 0

Column to use for index. If a sequence is given, a MultiIndex is used. Different default from `read_table`

**parse\_dates** : boolean, default True

Parse dates. Different default from read\_table

**Returns** y : DataFrame

#### Notes

Preferable to use read\_table for most general purposes but from\_csv makes for an easy roundtrip to and from file, especially with a DataFrame of time series data

### pandas.DataFrame.from\_records

**classmethod** DataFrame.**from\_records**(data, index=None, exclude=None, columns=None, names=None, coerce\_float=False)

Convert structured or record ndarray to DataFrame

**Parameters** data : ndarray (structured dtype), list of tuples, or DataFrame

index : string, list of fields, array-like

Field of array to use as the index, alternately a specific set of input labels to use

**exclude: sequence, default None :**

Columns or fields to exclude

**columns** : sequence, default None

Column names to use, replacing any found in passed data

**coerce\_float** : boolean, default False

Attempt to convert values to non-string, non-numeric objects (like decimal.Decimal) to floating point, useful for SQL result sets

**Returns** df : DataFrame

### pandas.DataFrame.to\_csv

DataFrame.**to\_csv**(path\_or\_buf, sep=',', na\_rep='', cols=None, header=True, index=True, index\_label=None, mode='w', nanRep=None, encoding=None)

Write DataFrame to a comma-separated values (csv) file

**Parameters** path\_or\_buf : string or file handle / StringIO

File path

**na\_rep** : string, default ''

Missing data representation

**cols** : sequence, optional

Columns to write

**header** : boolean or list of string, default True

Write out column names. If a list of string is given it is assumed to be aliases for the column names

**index** : boolean, default True

Write row names (index)

**index\_label** : string or sequence, default None

Column label for index column(s) if desired. If None is given, and *header* and *index* are True, then the index names are used. A sequence should be given if the DataFrame uses MultiIndex.

**mode** : Python write mode, default 'w'

**sep** : character, default ','

Field delimiter for the output file.

**encoding** : string, optional

a string representing the encoding to use if the contents are non-ascii, for python versions prior to 3

### pandas.DataFrame.to\_excel

`DataFrame.to_excel(excel_writer, sheet_name='sheet1', na_rep='', cols=None, header=True, index=True, index_label=None)`

Write DataFrame to a excel sheet

**Parameters** **excel\_writer** : string or ExcelWriter object

File path or existing ExcelWriter

**sheet\_name** : string, default 'sheet1'

Name of sheet which will contain DataFrame

**na\_rep** : string, default ''

Missing data rep'n

**cols** : sequence, optional

Columns to write

**header** : boolean or list of string, default True

Write out column names. If a list of string is given it is assumed to be aliases for the column names

**index** : boolean, default True

Write row names (index)

**index\_label** : string or sequence, default None

Column label for index column(s) if desired. If None is given, and *header* and *index* are True, then the index names are used. A sequence should be given if the DataFrame uses MultiIndex.

### Notes

If passing an existing ExcelWriter object, then the sheet will be added to the existing workbook. This can be used to save different DataFrames to one workbook >>> `writer = ExcelWriter('output.xlsx')` >>> `df1.to_excel(writer,'sheet1')` >>> `df2.to_excel(writer,'sheet2')` >>> `writer.save()`

### pandas.DataFrame.to\_dict

`DataFrame.to_dict (outtype='dict')`

Convert DataFrame to dictionary.

**Parameters** `outtype` : str {'dict', 'list', 'series'}

Determines the type of the values of the dictionary. The default *dict* is a nested dictionary {column -> {index -> value}}. *list* returns {column -> list(values)}. *series* returns {column -> Series(values)}. Abbreviations are allowed.

**Returns** `result` : dict like {column -> {index -> value}}

### pandas.DataFrame.to\_records

`DataFrame.to_records (index=True)`

Convert DataFrame to record array. Index will be put in the 'index' field of the record array if requested

**Parameters** `index` : boolean, default True

Include index in resulting record array, stored in 'index' field

**Returns** `y` : recarray

### pandas.DataFrame.to\_sparse

`DataFrame.to_sparse (fill_value=None, kind='block')`

Convert to SparseDataFrame

**Parameters** `fill_value` : float, default NaN

`kind` : {'block', 'integer'}

**Returns** `y` : SparseDataFrame

### pandas.DataFrame.to\_string

`DataFrame.to_string (buf=None, columns=None, col_space=None, colSpace=None, header=True, index=True, na_rep='NaN', formatters=None, float_format=None, sparsify=None, nanRep=None, index_names=True, justify=None, force_unicode=False)`

Render a DataFrame to a console-friendly tabular output.

**Parameters** `frame` : DataFrame

object to render

`buf` : StringIO-like, optional

buffer to write to

`columns` : sequence, optional

the subset of columns to write; default None writes all columns

`col_space` : int, optional

the width of each columns

`header` : bool, optional

whether to print column labels, default True

**index** : bool, optional

whether to print index (row) labels, default True

**na\_rep** : string, optional

string representation of NAN to use, default 'NaN'

**formatters** : list or dict of one-parameter functions, optional

formatter functions to apply to columns' elements by position or name, default None

**float\_format** : one-parameter function, optional

formatter function to apply to columns' elements if they are floats default None

**sparsify** : bool, optional

Set to False for a DataFrame with a hierarchical index to print every multiindex key at each row, default True

**justify** : { 'left', 'right' }, default None

Left or right-justify the column labels. If None uses the option from the configuration in pandas.core.common, 'left' out of the box

**index\_names** : bool, optional

Prints the names of the indexes, default True

**force\_unicode** : bool, default False

Always return a unicode result

**Returns** **formatted** : string (or unicode, depending on data and options)

## **pandas.DataFrame.save**

`DataFrame.save` (*path*)

## **pandas.DataFrame.load**

**classmethod** `DataFrame.load` (*path*)

## **pandas.DataFrame.info**

`DataFrame.info` (*verbose=True, buf=None*)

Concise summary of a DataFrame, used in `__repr__` when very large.

**Parameters** **verbose** : boolean, default True

If False, don't print column count summary

**buf** : writable buffer, defaults to `sys.stdout`

## 21.4 Panel

### 21.4.1 Computations / Descriptive Stats

# PYTHON MODULE INDEX

p

pandas, [1](#)





# PYTHON MODULE INDEX

p

pandas, [1](#)



# INDEX

## Symbols

`__init__()` (pandas.DataFrame method), 309  
`__init__()` (pandas.Series method), 287  
`__iter__()` (pandas.DataFrame method), 311  
`__iter__()` (pandas.Series method), 288

## A

`add()` (pandas.DataFrame method), 312  
`add()` (pandas.Series method), 288  
`add_prefix()` (pandas.DataFrame method), 325  
`add_suffix()` (pandas.DataFrame method), 325  
`align()` (pandas.DataFrame method), 325  
`align()` (pandas.Series method), 297  
`append()` (pandas.DataFrame method), 336  
`append()` (pandas.Series method), 303  
`apply()` (pandas.DataFrame method), 317  
`apply()` (pandas.Series method), 290  
`applymap()` (pandas.DataFrame method), 317  
`argsort()` (pandas.Series method), 301  
`as_matrix()` (pandas.DataFrame method), 308  
`asfreq()` (pandas.DataFrame method), 336  
`asfreq()` (pandas.Series method), 303  
`asof()` (pandas.Series method), 304  
`astype()` (pandas.DataFrame method), 310  
`astype()` (pandas.Series method), 287  
`autocorr()` (pandas.Series method), 292  
`axes` (pandas.DataFrame attribute), 309

## C

`clip()` (pandas.DataFrame method), 319  
`clip()` (pandas.Series method), 292  
`clip_lower()` (pandas.DataFrame method), 319  
`clip_lower()` (pandas.Series method), 292  
`clip_upper()` (pandas.DataFrame method), 319  
`clip_upper()` (pandas.Series method), 293  
`combine()` (pandas.DataFrame method), 316  
`combine()` (pandas.Series method), 289  
`combine_first()` (pandas.DataFrame method), 316  
`combine_first()` (pandas.Series method), 290  
`combineAdd()` (pandas.DataFrame method), 316  
`combineMult()` (pandas.DataFrame method), 316

`concat()` (in module pandas.tools.merge), 271  
`copy()` (pandas.DataFrame method), 310  
`copy()` (pandas.Series method), 287  
`corr()` (pandas.DataFrame method), 319  
`corr()` (pandas.Series method), 293  
`corrwith()` (pandas.DataFrame method), 319  
`count()` (pandas.DataFrame method), 320  
`count()` (pandas.Series method), 293  
`cumprod()` (pandas.DataFrame method), 320  
`cumprod()` (pandas.Series method), 293  
`cumsum()` (pandas.DataFrame method), 320  
`cumsum()` (pandas.Series method), 293

## D

`delevel()` (pandas.DataFrame method), 331  
`describe()` (pandas.DataFrame method), 321  
`describe()` (pandas.Series method), 294  
`diff()` (pandas.DataFrame method), 321  
`diff()` (pandas.Series method), 294  
`div()` (pandas.DataFrame method), 312  
`div()` (pandas.Series method), 289  
`drop()` (pandas.DataFrame method), 326  
`drop()` (pandas.Series method), 298  
`dropna()` (pandas.DataFrame method), 329  
`dropna()` (pandas.Series method), 300  
`dtype` (pandas.Series attribute), 286  
`dtypes` (pandas.DataFrame attribute), 308

## E

`ewma()` (in module pandas.stats.moments), 282  
`ewmcorr()` (in module pandas.stats.moments), 284  
`ewmcov()` (in module pandas.stats.moments), 285  
`ewmstd()` (in module pandas.stats.moments), 283  
`ewmvar()` (in module pandas.stats.moments), 283

## F

`fillna()` (pandas.DataFrame method), 330  
`fillna()` (pandas.Series method), 300  
`filter()` (pandas.DataFrame method), 326  
`first_valid_index()` (pandas.DataFrame method), 337  
`first_valid_index()` (pandas.Series method), 304

from\_csv() (pandas.DataFrame class method), 339  
from\_csv() (pandas.Series class method), 306  
from\_records() (pandas.DataFrame class method), 340

## G

get() (pandas.io.pytables.HDFStore method), 277  
get() (pandas.Series method), 288  
get\_dtype\_counts() (pandas.DataFrame method), 309  
groupby() (pandas.DataFrame method), 318  
groupby() (pandas.Series method), 291

## H

head() (pandas.DataFrame method), 329  
hist() (pandas.DataFrame method), 337  
hist() (pandas.Series method), 305

## I

info() (pandas.DataFrame method), 343  
insert() (pandas.DataFrame method), 311  
interpolate() (pandas.Series method), 301  
isnull() (pandas.Series method), 286  
iteritems() (pandas.DataFrame method), 311  
iteritems() (pandas.Series method), 288  
ix (pandas.DataFrame attribute), 311  
ix (pandas.Series attribute), 288

## J

join() (pandas.DataFrame method), 334

## L

last\_valid\_index() (pandas.DataFrame method), 337  
last\_valid\_index() (pandas.Series method), 304  
load() (in module pandas.core.common), 272  
load() (pandas.DataFrame class method), 343  
load() (pandas.Series class method), 307

## M

mad() (pandas.DataFrame method), 321  
map() (pandas.Series method), 290  
max() (pandas.DataFrame method), 321  
max() (pandas.Series method), 294  
mean() (pandas.DataFrame method), 322  
mean() (pandas.Series method), 294  
median() (pandas.DataFrame method), 322  
median() (pandas.Series method), 295  
merge() (in module pandas.tools.merge), 270  
merge() (pandas.DataFrame method), 334  
min() (pandas.DataFrame method), 322  
min() (pandas.Series method), 295  
mul() (pandas.DataFrame method), 313  
mul() (pandas.Series method), 289

## N

ndim (pandas.DataFrame attribute), 309

notnull() (pandas.Series method), 286

## O

order() (pandas.Series method), 301

## P

pandas (module), 1  
parse() (pandas.io.parsers.ExcelFile method), 276  
pivot() (pandas.DataFrame method), 331  
pivot\_table() (in module pandas.tools.pivot), 269  
plot() (pandas.DataFrame method), 338  
plot() (pandas.Series method), 305  
pop() (pandas.DataFrame method), 311  
prod() (pandas.DataFrame method), 323  
prod() (pandas.Series method), 295  
put() (pandas.io.pytables.HDFStore method), 277

## Q

quantile() (pandas.DataFrame method), 323  
quantile() (pandas.Series method), 295

## R

radd() (pandas.DataFrame method), 314  
rdiv() (pandas.DataFrame method), 314  
read\_csv() (in module pandas.io.parsers), 275  
read\_table() (in module pandas.io.parsers), 273  
reindex() (pandas.DataFrame method), 327  
reindex() (pandas.Series method), 298  
reindex\_like() (pandas.DataFrame method), 327  
reindex\_like() (pandas.Series method), 298  
rename() (pandas.DataFrame method), 328  
rename() (pandas.Series method), 299  
rmul() (pandas.DataFrame method), 315  
rolling\_apply() (in module pandas.stats.moments), 281  
rolling\_corr() (in module pandas.stats.moments), 280  
rolling\_count() (in module pandas.stats.moments), 278  
rolling\_cov() (in module pandas.stats.moments), 280  
rolling\_kurt() (in module pandas.stats.moments), 281  
rolling\_mean() (in module pandas.stats.moments), 278  
rolling\_median() (in module pandas.stats.moments), 279  
rolling\_quantile() (in module pandas.stats.moments), 281  
rolling\_skew() (in module pandas.stats.moments), 280  
rolling\_std() (in module pandas.stats.moments), 279  
rolling\_sum() (in module pandas.stats.moments), 278  
rolling\_var() (in module pandas.stats.moments), 279  
rsub() (pandas.DataFrame method), 315

## S

save() (in module pandas.core.common), 273  
save() (pandas.DataFrame method), 343  
save() (pandas.Series method), 307  
select() (pandas.DataFrame method), 328  
select() (pandas.Series method), 299

shape (pandas.DataFrame attribute), 309  
shift() (pandas.DataFrame method), 336  
shift() (pandas.Series method), 304  
skew() (pandas.DataFrame method), 323  
skew() (pandas.Series method), 296  
sort() (pandas.Series method), 302  
sort\_index() (pandas.DataFrame method), 330  
sort\_index() (pandas.Series method), 302  
sortlevel() (pandas.DataFrame method), 332  
sortlevel() (pandas.Series method), 302  
stack() (pandas.DataFrame method), 332  
std() (pandas.DataFrame method), 324  
std() (pandas.Series method), 296  
sub() (pandas.DataFrame method), 313  
sub() (pandas.Series method), 289  
sum() (pandas.DataFrame method), 323  
sum() (pandas.Series method), 296  
swaplevel() (pandas.DataFrame method), 332

## T

T (pandas.DataFrame attribute), 333  
tail() (pandas.DataFrame method), 329  
take() (pandas.DataFrame method), 328  
take() (pandas.Series method), 299  
to\_csv() (pandas.DataFrame method), 340  
to\_csv() (pandas.Series method), 307  
to\_dict() (pandas.DataFrame method), 342  
to\_dict() (pandas.Series method), 308  
to\_excel() (pandas.DataFrame method), 341  
to\_records() (pandas.DataFrame method), 342  
to\_sparse() (pandas.DataFrame method), 342  
to\_sparse() (pandas.Series method), 308  
to\_string() (pandas.DataFrame method), 342  
transpose() (pandas.DataFrame method), 333  
truncate() (pandas.DataFrame method), 329  
truncate() (pandas.Series method), 300

## U

unstack() (pandas.DataFrame method), 333  
unstack() (pandas.Series method), 302

## V

value\_counts() (pandas.Series method), 297  
values (pandas.DataFrame attribute), 309  
values (pandas.Series attribute), 286  
var() (pandas.DataFrame method), 324  
var() (pandas.Series method), 296

## W

weekday (pandas.Series attribute), 304

## X

xs() (pandas.DataFrame method), 311