```
(-0.523, 0.0296]   -0.260266
(0.0296, 0.654]     0.361802
(0.654, 2.21]       1.073801
dtype: float64
```

### Grouping with a grouper specification

You may need to specify a bit more data to properly group. You can use the `pd.Grouper` to provide this local control.

```
In [171]: import datetime

In [172]: df = pd.DataFrame({'Branch': 'A A A A A A A B'.split(),
   .....:                    'Buyer': 'Carl Mark Carl Carl Joe Joe Joe Carl'.split(),
   .....:                    'Quantity': [1, 3, 5, 1, 8, 1, 9, 3],
   .....:                    'Date': [
   .....:                        datetime.datetime(2013, 1, 1, 13, 0),
   .....:                        datetime.datetime(2013, 1, 1, 13, 5),
   .....:                        datetime.datetime(2013, 10, 1, 20, 0),
   .....:                        datetime.datetime(2013, 10, 2, 10, 0),
   .....:                        datetime.datetime(2013, 10, 1, 20, 0),
   .....:                        datetime.datetime(2013, 10, 2, 10, 0),
   .....:                        datetime.datetime(2013, 12, 2, 12, 0),
   .....:                        datetime.datetime(2013, 12, 2, 14, 0)]
   .....:                   })
   .....:

In [173]: df
Out[173]:
  Branch Buyer  Quantity                Date
0      A  Carl         1 2013-01-01 13:00:00
1      A  Mark         3 2013-01-01 13:05:00
2      A  Carl         5 2013-10-01 20:00:00
3      A  Carl         1 2013-10-02 10:00:00
4      A   Joe         8 2013-10-01 20:00:00
5      A   Joe         1 2013-10-02 10:00:00
6      A   Joe         9 2013-12-02 12:00:00
7      B  Carl         3 2013-12-02 14:00:00
```

Groupby a specific column with the desired frequency. This is like resampling.

```
In [174]: df.groupby([pd.Grouper(freq='1M', key='Date'), 'Buyer']).sum()
Out[174]:
                 Quantity
Date       Buyer
2013-01-31 Carl         1
           Mark         3
2013-10-31 Carl         6
           Joe          9
2013-12-31 Carl         3
           Joe          9
```

You have an ambiguous specification in that you have a named index and a column that could be potential groupers.

```
In [175]: df = df.set_index('Date')

In [176]: df['Date'] = df.index + pd.offsets.MonthEnd(2)

In [177]: df.groupby([pd.Grouper(freq='6M', key='Date'), 'Buyer']).sum()
Out[177]:
                  Quantity
Date       Buyer
2013-02-28 Carl          1
           Mark          3
2014-02-28 Carl          9
           Joe          18

In [178]: df.groupby([pd.Grouper(freq='6M', level='Date'), 'Buyer']).sum()
Out[178]:
                  Quantity
Date       Buyer
2013-01-31 Carl          1
           Mark          3
2014-01-31 Carl          9
           Joe          18
```

**Taking the first rows of each group**

Just like for a DataFrame or Series you can call head and tail on a groupby:

```
In [179]: df = pd.DataFrame([[1, 2], [1, 4], [5, 6]], columns=['A', 'B'])

In [180]: df
Out[180]:
   A  B
0  1  2
1  1  4
2  5  6

In [181]: g = df.groupby('A')

In [182]: g.head(1)
Out[182]:
   A  B
0  1  2
2  5  6

In [183]: g.tail(1)
Out[183]:
   A  B
1  1  4
2  5  6
```

This shows the first or last n rows from each group.

### Taking the nth row of each group

To select from a DataFrame or Series the nth item, use `nth()`. This is a reduction method, and will return a single row (or no row) per group if you pass an int for n:

```
In [184]: df = pd.DataFrame([[1, np.nan], [1, 4], [5, 6]], columns=['A', 'B'])

In [185]: g = df.groupby('A')

In [186]: g.nth(0)
Out[186]:
     B
A
1  NaN
5  6.0

In [187]: g.nth(-1)
Out[187]:
     B
A
1  4.0
5  6.0

In [188]: g.nth(1)
Out[188]:
     B
A
1  4.0
```

If you want to select the nth not-null item, use the `dropna` kwarg. For a DataFrame this should be either `'any'` or `'all'` just like you would pass to dropna:

```
# nth(0) is the same as g.first()
In [189]: g.nth(0, dropna='any')
Out[189]:
     B
A
1  4.0
5  6.0

In [190]: g.first()
Out[190]:
     B
A
1  4.0
5  6.0

# nth(-1) is the same as g.last()
In [191]: g.nth(-1, dropna='any')  # NaNs denote group exhausted when using dropna
Out[191]:
     B
A
1  4.0
5  6.0

In [192]: g.last()
Out[192]:
     B
```

(continues on next page)

```
A
1  4.0
5  6.0

In [193]: g.B.nth(0, dropna='all')
Out[193]:
A
1    4.0
5    6.0
Name: B, dtype: float64
```

As with other methods, passing `as_index=False`, will achieve a filtration, which returns the grouped row.

```
In [194]: df = pd.DataFrame([[1, np.nan], [1, 4], [5, 6]], columns=['A', 'B'])

In [195]: g = df.groupby('A', as_index=False)

In [196]: g.nth(0)
Out[196]:
   A    B
0  1  NaN
2  5  6.0

In [197]: g.nth(-1)
Out[197]:
   A    B
1  1  4.0
2  5  6.0
```

You can also select multiple rows from each group by specifying multiple nth values as a list of ints.

```
In [198]: business_dates = pd.date_range(start='4/1/2014', end='6/30/2014', freq='B')

In [199]: df = pd.DataFrame(1, index=business_dates, columns=['a', 'b'])

# get the first, 4th, and last date index for each month
In [200]: df.groupby([df.index.year, df.index.month]).nth([0, 3, -1])
Out[200]:
        a  b
2014 4  1  1
     4  1  1
     4  1  1
     5  1  1
     5  1  1
     5  1  1
     6  1  1
     6  1  1
     6  1  1
```

### Enumerate group items

To see the order in which each row appears within its group, use the `cumcount` method:

```
In [201]: dfg = pd.DataFrame(list('aaabba'), columns=['A'])

In [202]: dfg
Out[202]:
   A
0  a
1  a
2  a
3  b
4  b
5  a

In [203]: dfg.groupby('A').cumcount()
Out[203]:
0    0
1    1
2    2
3    0
4    1
5    3
dtype: int64

In [204]: dfg.groupby('A').cumcount(ascending=False)
Out[204]:
0    3
1    2
2    1
3    1
4    0
5    0
dtype: int64
```

### Enumerate groups

To see the ordering of the groups (as opposed to the order of rows within a group given by `cumcount`) you can use `ngroup()`.

Note that the numbers given to the groups match the order in which the groups would be seen when iterating over the groupby object, not the order they are first observed.

```
In [205]: dfg = pd.DataFrame(list('aaabba'), columns=['A'])

In [206]: dfg
Out[206]:
   A
0  a
1  a
2  a
3  b
4  b
5  a
```

```
In [207]: dfg.groupby('A').ngroup()
Out[207]:
0    0
1    0
2    0
3    1
4    1
5    0
dtype: int64

In [208]: dfg.groupby('A').ngroup(ascending=False)
Out[208]:
0    1
1    1
2    1
3    0
4    0
5    1
dtype: int64
```

### Plotting

Groupby also works with some plotting methods. For example, suppose we suspect that some features in a DataFrame may differ by group, in this case, the values in column 1 where the group is "B" are 3 higher on average.

```
In [209]: np.random.seed(1234)

In [210]: df = pd.DataFrame(np.random.randn(50, 2))

In [211]: df['g'] = np.random.choice(['A', 'B'], size=50)

In [212]: df.loc[df['g'] == 'B', 1] += 3
```
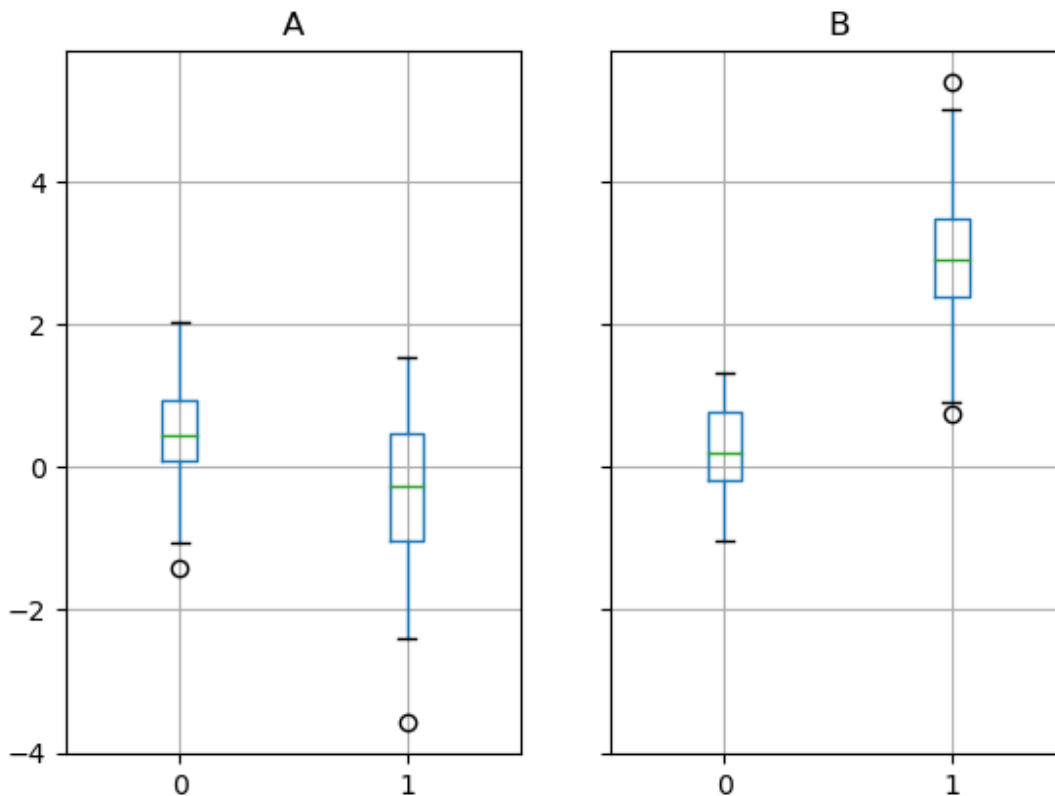
We can easily visualize this with a boxplot:

```
In [213]: df.groupby('g').boxplot()
Out[213]:
A       AxesSubplot(0.1,0.15;0.363636x0.75)
B    AxesSubplot(0.536364,0.15;0.363636x0.75)
dtype: object
```

The result of calling `boxplot` is a dictionary whose keys are the values of our grouping column `g` ("A" and "B"). The values of the resulting dictionary can be controlled by the `return_type` keyword of `boxplot`. See the *visualization documentation* for more.

> **Warning:** For historical reasons, `df.groupby("g").boxplot()` is not equivalent to `df.boxplot(by= "g")`. See *here* for an explanation.

### Piping function calls

New in version 0.21.0.

Similar to the functionality provided by `DataFrame` and `Series`, functions that take `GroupBy` objects can be chained together using a `pipe` method to allow for a cleaner, more readable syntax. To read about `.pipe` in general terms, see *here*.

Combining `.groupby` and `.pipe` is often useful when you need to reuse GroupBy objects.

As an example, imagine having a DataFrame with columns for stores, products, revenue and quantity sold. We'd like to do a groupwise calculation of *prices* (i.e. revenue/quantity) per store and per product. We could do this in a multi-step operation, but expressing it in terms of piping can make the code more readable. First we set the data:

```
In [214]: n = 1000
```

(continues on next page)

```
In [215]: df = pd.DataFrame({'Store': np.random.choice(['Store_1', 'Store_2'], n),
   .....:                      'Product': np.random.choice(['Product_1',
   .....:                                                   'Product_2'], n),
   .....:                      'Revenue': (np.random.random(n) * 50 + 10).round(2),
   .....:                      'Quantity': np.random.randint(1, 10, size=n)})
   .....:

In [216]: df.head(2)
Out[216]:
    Store    Product   Revenue  Quantity
0  Store_2  Product_1    26.12         1
1  Store_2  Product_1    28.86         1
```

Now, to find prices per store/product, we can simply do:

```
In [217]: (df.groupby(['Store', 'Product'])
   .....:      .pipe(lambda grp: grp.Revenue.sum() / grp.Quantity.sum())
   .....:      .unstack().round(2))
   .....:
Out[217]:
Product  Product_1  Product_2
Store
Store_1       6.82       7.05
Store_2       6.30       6.64
```

Piping can also be expressive when you want to deliver a grouped object to some arbitrary function, for example:

```
In [218]: def mean(groupby):
   .....:       return groupby.mean()
   .....:

In [219]: df.groupby(['Store', 'Product']).pipe(mean)
Out[219]:
                    Revenue  Quantity
Store    Product
Store_1  Product_1  34.622727  5.075758
         Product_2  35.482815  5.029630
Store_2  Product_1  32.972837  5.237589
         Product_2  34.684360  5.224000
```

where `mean` takes a GroupBy object and finds the mean of the Revenue and Quantity columns respectively for each Store-Product combination. The `mean` function can be any function that takes in a GroupBy object; the `.pipe` will pass the GroupBy object as a parameter into the function you specify.

### 2.13.10 Examples

**Regrouping by factor**

Regroup columns of a DataFrame according to their sum, and sum the aggregated ones.

```
In [220]: df = pd.DataFrame({'a': [1, 0, 0], 'b': [0, 1, 0],
   .....:                      'c': [1, 0, 0], 'd': [2, 3, 4]})
   .....:

In [221]: df
```

```
Out[221]:
   a  b  c  d
0  1  0  1  2
1  0  1  0  3
2  0  0  0  4

In [222]: df.groupby(df.sum(), axis=1).sum()
Out[222]:
   1  9
0  2  2
1  1  3
2  0  4
```

### Multi-column factorization

By using `ngroup()`, we can extract information about the groups in a way similar to *factorize()* (as described further in the *reshaping API*) but which applies naturally to multiple columns of mixed type and different sources. This can be useful as an intermediate categorical-like step in processing, when the relationships between the group rows are more important than their content, or as input to an algorithm which only accepts the integer encoding. (For more information about support in pandas for full categorical data, see the *Categorical introduction* and the *API documentation*.)

```
In [223]: dfg = pd.DataFrame({"A": [1, 1, 2, 3, 2], "B": list("aaaba")})

In [224]: dfg
Out[224]:
   A  B
0  1  a
1  1  a
2  2  a
3  3  b
4  2  a

In [225]: dfg.groupby(["A", "B"]).ngroup()
Out[225]:
0    0
1    0
2    1
3    2
4    1
dtype: int64

In [226]: dfg.groupby(["A", [0, 0, 0, 1, 1]]).ngroup()
Out[226]:
0    0
1    0
2    1
3    3
4    2
dtype: int64
```

### Groupby by indexer to 'resample' data

Resampling produces new hypothetical samples (resamples) from already existing observed data or from a model that generates data. These new samples are similar to the pre-existing samples.

In order to resample to work on indices that are non-datetimelike, the following procedure can be utilized.

In the following examples, **df.index // 5** returns a binary array which is used to determine what gets selected for the groupby operation.

---

**Note:** The below example shows how we can downsample by consolidation of samples into fewer samples. Here by using **df.index // 5**, we are aggregating the samples in bins. By applying **std()** function, we aggregate the information contained in many samples into a small subset of values which is their standard deviation thereby reducing the number of samples.

---

```
In [227]: df = pd.DataFrame(np.random.randn(10, 2))

In [228]: df
Out[228]:
          0         1
0 -0.793893  0.321153
1  0.342250  1.618906
2 -0.975807  1.918201
3 -0.810847 -1.405919
4 -1.977759  0.461659
5  0.730057 -1.316938
6 -0.751328  0.528290
7 -0.257759 -1.081009
8  0.505895 -1.701948
9 -1.006349  0.020208

In [229]: df.index // 5
Out[229]: Int64Index([0, 0, 0, 0, 0, 1, 1, 1, 1, 1], dtype='int64')

In [230]: df.groupby(df.index // 5).std()
Out[230]:
          0         1
0  0.823647  1.312912
1  0.760109  0.942941
```

### Returning a Series to propagate names

Group DataFrame columns, compute a set of metrics and return a named Series. The Series name is used as the name for the column index. This is especially useful in conjunction with reshaping operations such as stacking in which the column index name will be used as the name of the inserted column:

```
In [231]: df = pd.DataFrame({'a': [0, 0, 0, 0, 1, 1, 1, 1, 2, 2, 2, 2],
   .....:                    'b': [0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1],
   .....:                    'c': [1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0],
   .....:                    'd': [0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1]})
   .....:

In [232]: def compute_metrics(x):
   .....:     result = {'b_sum': x['b'].sum(), 'c_mean': x['c'].mean()}
   .....:     return pd.Series(result, name='metrics')
```

```
   .....:

In [233]: result = df.groupby('a').apply(compute_metrics)

In [234]: result
Out[234]:
metrics  b_sum  c_mean
a
0          2.0     0.5
1          2.0     0.5
2          2.0     0.5

In [235]: result.stack()
Out[235]:
a  metrics
0  b_sum     2.0
   c_mean    0.5
1  b_sum     2.0
   c_mean    0.5
2  b_sum     2.0
   c_mean    0.5
dtype: float64
```

## 2.14 Time series / date functionality

pandas contains extensive capabilities and features for working with time series data for all domains. Using the NumPy `datetime64` and `timedelta64` dtypes, pandas has 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.

For example, pandas supports:

Parsing time series information from various sources and formats

```
In [1]: import datetime

In [2]: dti = pd.to_datetime(['1/1/2018', np.datetime64('2018-01-01'),
   ...:                        datetime.datetime(2018, 1, 1)])
   ...:

In [3]: dti
Out[3]: DatetimeIndex(['2018-01-01', '2018-01-01', '2018-01-01'], dtype=
→'datetime64[ns]', freq=None)
```

Generate sequences of fixed-frequency dates and time spans

```
In [4]: dti = pd.date_range('2018-01-01', periods=3, freq='H')

In [5]: dti
Out[5]:
DatetimeIndex(['2018-01-01 00:00:00', '2018-01-01 01:00:00',
               '2018-01-01 02:00:00'],
              dtype='datetime64[ns]', freq='H')
```

Manipulating and converting date times with timezone information

---

```
In [6]: dti = dti.tz_localize('UTC')

In [7]: dti
Out[7]:
DatetimeIndex(['2018-01-01 00:00:00+00:00', '2018-01-01 01:00:00+00:00',
               '2018-01-01 02:00:00+00:00'],
              dtype='datetime64[ns, UTC]', freq='H')

In [8]: dti.tz_convert('US/Pacific')
Out[8]:
DatetimeIndex(['2017-12-31 16:00:00-08:00', '2017-12-31 17:00:00-08:00',
               '2017-12-31 18:00:00-08:00'],
              dtype='datetime64[ns, US/Pacific]', freq='H')
```

Resampling or converting a time series to a particular frequency

```
In [9]: idx = pd.date_range('2018-01-01', periods=5, freq='H')

In [10]: ts = pd.Series(range(len(idx)), index=idx)

In [11]: ts
Out[11]:
2018-01-01 00:00:00    0
2018-01-01 01:00:00    1
2018-01-01 02:00:00    2
2018-01-01 03:00:00    3
2018-01-01 04:00:00    4
Freq: H, dtype: int64

In [12]: ts.resample('2H').mean()
Out[12]:
2018-01-01 00:00:00    0.5
2018-01-01 02:00:00    2.5
2018-01-01 04:00:00    4.0
Freq: 2H, dtype: float64
```

Performing date and time arithmetic with absolute or relative time increments

```
In [13]: friday = pd.Timestamp('2018-01-05')

In [14]: friday.day_name()
Out[14]: 'Friday'

# Add 1 day
In [15]: saturday = friday + pd.Timedelta('1 day')

In [16]: saturday.day_name()
Out[16]: 'Saturday'

# Add 1 business day (Friday --> Monday)
In [17]: monday = friday + pd.offsets.BDay()

In [18]: monday.day_name()
Out[18]: 'Monday'
```

pandas provides a relatively compact and self-contained set of tools for performing the above tasks and more.

### 2.14.1 Overview

pandas captures 4 general time related concepts:

1. Date times: A specific date and time with timezone support. Similar to `datetime.datetime` from the standard library.

2. Time deltas: An absolute time duration. Similar to `datetime.timedelta` from the standard library.

3. Time spans: A span of time defined by a point in time and its associated frequency.

4. Date offsets: A relative time duration that respects calendar arithmetic. Similar to `dateutil.relativedelta.relativedelta` from the `dateutil` package.

| Concept | Scalar Class | Array Class | pandas Data Type | Primary Creation Method |
|---------|--------------|-------------|------------------|-------------------------|
| Date times | `Timestamp` | `DatetimeIndex` | `datetime64[ns]` or `datetime64[ns, tz]` | `to_datetime` or `date_range` |
| Time deltas | `Timedelta` | `TimedeltaIndex` | `timedelta64[ns]` | `to_timedelta` or `timedelta_range` |
| Time spans | `Period` | `PeriodIndex` | `period[freq]` | `Period` or `period_range` |
| Date offsets | `DateOffset` | `None` | `None` | `DateOffset` |

For time series data, it's conventional to represent the time component in the index of a *Series* or *DataFrame* so manipulations can be performed with respect to the time element.

```
In [19]: pd.Series(range(3), index=pd.date_range('2000', freq='D', periods=3))
Out[19]:
2000-01-01    0
2000-01-02    1
2000-01-03    2
Freq: D, dtype: int64
```

However, *Series* and *DataFrame* can directly also support the time component as data itself.

```
In [20]: pd.Series(pd.date_range('2000', freq='D', periods=3))
Out[20]:
0   2000-01-01
1   2000-01-02
2   2000-01-03
dtype: datetime64[ns]
```

*Series* and *DataFrame* have extended data type support and functionality for `datetime`, `timedelta` and `Period` data when passed into those constructors. `DateOffset` data however will be stored as `object` data.

```
In [21]: pd.Series(pd.period_range('1/1/2011', freq='M', periods=3))
Out[21]:
0    2011-01
1    2011-02
2    2011-03
dtype: period[M]

In [22]: pd.Series([pd.DateOffset(1), pd.DateOffset(2)])
Out[22]:
0        <DateOffset>
```

(continues on next page)

```
1    <2 * DateOffsets>
dtype: object

In [23]: pd.Series(pd.date_range('1/1/2011', freq='M', periods=3))
Out[23]:
0   2011-01-31
1   2011-02-28
2   2011-03-31
dtype: datetime64[ns]
```

Lastly, pandas represents null date times, time deltas, and time spans as `NaT` which is useful for representing missing or null date like values and behaves similar as `np.nan` does for float data.

```
In [24]: pd.Timestamp(pd.NaT)
Out[24]: NaT

In [25]: pd.Timedelta(pd.NaT)
Out[25]: NaT

In [26]: pd.Period(pd.NaT)
Out[26]: NaT

# Equality acts as np.nan would
In [27]: pd.NaT == pd.NaT
Out[27]: False
```

## 2.14.2 Timestamps vs. Time Spans

Timestamped data is the most basic type of time series data that associates values with points in time. For pandas objects it means using the points in time.

```
In [28]: pd.Timestamp(datetime.datetime(2012, 5, 1))
Out[28]: Timestamp('2012-05-01 00:00:00')

In [29]: pd.Timestamp('2012-05-01')
Out[29]: Timestamp('2012-05-01 00:00:00')

In [30]: pd.Timestamp(2012, 5, 1)
Out[30]: Timestamp('2012-05-01 00:00:00')
```

However, in many cases it is more natural to associate things like change variables with a time span instead. The span represented by `Period` can be specified explicitly, or inferred from datetime string format.

For example:

```
In [31]: pd.Period('2011-01')
Out[31]: Period('2011-01', 'M')

In [32]: pd.Period('2012-05', freq='D')
Out[32]: Period('2012-05-01', 'D')
```

*Timestamp* and *Period* can serve as an index. Lists of `Timestamp` and `Period` are automatically coerced to *DatetimeIndex* and *PeriodIndex* respectively.

```
In [33]: dates = [pd.Timestamp('2012-05-01'),
   ....:          pd.Timestamp('2012-05-02'),
   ....:          pd.Timestamp('2012-05-03')]
   ....:

In [34]: ts = pd.Series(np.random.randn(3), dates)

In [35]: type(ts.index)
Out[35]: pandas.core.indexes.datetimes.DatetimeIndex

In [36]: ts.index
Out[36]: DatetimeIndex(['2012-05-01', '2012-05-02', '2012-05-03'], dtype=
→'datetime64[ns]', freq=None)

In [37]: ts
Out[37]:
2012-05-01    0.469112
2012-05-02   -0.282863
2012-05-03   -1.509059
dtype: float64

In [38]: periods = [pd.Period('2012-01'), pd.Period('2012-02'), pd.Period('2012-03')]

In [39]: ts = pd.Series(np.random.randn(3), periods)

In [40]: type(ts.index)
Out[40]: pandas.core.indexes.period.PeriodIndex

In [41]: ts.index
Out[41]: PeriodIndex(['2012-01', '2012-02', '2012-03'], dtype='period[M]', freq='M')

In [42]: ts
Out[42]:
2012-01   -1.135632
2012-02    1.212112
2012-03   -0.173215
Freq: M, dtype: float64
```

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.

### 2.14.3 Converting to timestamps

To convert a *Series* or list-like object of date-like objects e.g. strings, epochs, or a mixture, you can use the `to_datetime` function. When passed a `Series`, this returns a `Series` (with the same index), while a list-like is converted to a `DatetimeIndex`:

```
In [43]: pd.to_datetime(pd.Series(['Jul 31, 2009', '2010-01-10', None]))
Out[43]:
0   2009-07-31
1   2010-01-10
2          NaT
dtype: datetime64[ns]
```

```
In [44]: pd.to_datetime(['2005/11/23', '2010.12.31'])
Out[44]: DatetimeIndex(['2005-11-23', '2010-12-31'], dtype='datetime64[ns]',
→freq=None)
```

If you use dates which start with the day first (i.e. European style), you can pass the `dayfirst` flag:

```
In [45]: pd.to_datetime(['04-01-2012 10:00'], dayfirst=True)
Out[45]: DatetimeIndex(['2012-01-04 10:00:00'], dtype='datetime64[ns]', freq=None)

In [46]: pd.to_datetime(['14-01-2012', '01-14-2012'], dayfirst=True)
Out[46]: DatetimeIndex(['2012-01-14', '2012-01-14'], dtype='datetime64[ns]',
→freq=None)
```

> **Warning:** You see in the above example that `dayfirst` isn't strict, so if a date can't be parsed with the day being first it will be parsed as if `dayfirst` were False.

If you pass a single string to `to_datetime`, it returns a single `Timestamp`. `Timestamp` can also accept string input, but it doesn't accept string parsing options like `dayfirst` or `format`, so use `to_datetime` if these are required.

```
In [47]: pd.to_datetime('2010/11/12')
Out[47]: Timestamp('2010-11-12 00:00:00')

In [48]: pd.Timestamp('2010/11/12')
Out[48]: Timestamp('2010-11-12 00:00:00')
```

You can also use the `DatetimeIndex` constructor directly:

```
In [49]: pd.DatetimeIndex(['2018-01-01', '2018-01-03', '2018-01-05'])
Out[49]: DatetimeIndex(['2018-01-01', '2018-01-03', '2018-01-05'], dtype=
→'datetime64[ns]', freq=None)
```

The string 'infer' can be passed in order to set the frequency of the index as the inferred frequency upon creation:

```
In [50]: pd.DatetimeIndex(['2018-01-01', '2018-01-03', '2018-01-05'], freq='infer')
Out[50]: DatetimeIndex(['2018-01-01', '2018-01-03', '2018-01-05'], dtype=
→'datetime64[ns]', freq='2D')
```

### Providing a format argument

In addition to the required datetime string, a `format` argument can be passed to ensure specific parsing. This could also potentially speed up the conversion considerably.

```
In [51]: pd.to_datetime('2010/11/12', format='%Y/%m/%d')
Out[51]: Timestamp('2010-11-12 00:00:00')

In [52]: pd.to_datetime('12-11-2010 00:00', format='%d-%m-%Y %H:%M')
Out[52]: Timestamp('2010-11-12 00:00:00')
```

For more information on the choices available when specifying the `format` option, see the Python datetime documentation.

---

### Assembling datetime from multiple DataFrame columns

You can also pass a `DataFrame` of integer or string columns to assemble into a `Series` of `Timestamps`.

```
In [53]: df = pd.DataFrame({'year': [2015, 2016],
   ....:                     'month': [2, 3],
   ....:                     'day': [4, 5],
   ....:                     'hour': [2, 3]})
   ....:

In [54]: pd.to_datetime(df)
Out[54]:
0   2015-02-04 02:00:00
1   2016-03-05 03:00:00
dtype: datetime64[ns]
```

You can pass only the columns that you need to assemble.

```
In [55]: pd.to_datetime(df[['year', 'month', 'day']])
Out[55]:
0   2015-02-04
1   2016-03-05
dtype: datetime64[ns]
```

`pd.to_datetime` looks for standard designations of the datetime component in the column names, including:

- required: `year`, `month`, `day`
- optional: `hour`, `minute`, `second`, `millisecond`, `microsecond`, `nanosecond`

### Invalid data

The default behavior, `errors='raise'`, is to raise when unparseable:

```
In [2]: pd.to_datetime(['2009/07/31', 'asd'], errors='raise')
ValueError: Unknown string format
```

Pass `errors='ignore'` to return the original input when unparseable:

```
In [56]: pd.to_datetime(['2009/07/31', 'asd'], errors='ignore')
Out[56]: Index(['2009/07/31', 'asd'], dtype='object')
```

Pass `errors='coerce'` to convert unparseable data to `NaT` (not a time):

```
In [57]: pd.to_datetime(['2009/07/31', 'asd'], errors='coerce')
Out[57]: DatetimeIndex(['2009-07-31', 'NaT'], dtype='datetime64[ns]', freq=None)
```

**Epoch timestamps**

pandas supports converting integer or float epoch times to `Timestamp` and `DatetimeIndex`. The default unit is nanoseconds, since that is how `Timestamp` objects are stored internally. However, epochs are often stored in another `unit` which can be specified. These are computed from the starting point specified by the `origin` parameter.

```
In [58]: pd.to_datetime([1349720105, 1349806505, 1349892905,
   ....:                 1349979305, 1350065705], unit='s')
   ....:
Out[58]:
DatetimeIndex(['2012-10-08 18:15:05', '2012-10-09 18:15:05',
               '2012-10-10 18:15:05', '2012-10-11 18:15:05',
               '2012-10-12 18:15:05'],
              dtype='datetime64[ns]', freq=None)

In [59]: pd.to_datetime([1349720105100, 1349720105200, 1349720105300,
   ....:                 1349720105400, 1349720105500], unit='ms')
   ....:
Out[59]:
DatetimeIndex(['2012-10-08 18:15:05.100000', '2012-10-08 18:15:05.200000',
               '2012-10-08 18:15:05.300000', '2012-10-08 18:15:05.400000',
               '2012-10-08 18:15:05.500000'],
              dtype='datetime64[ns]', freq=None)
```

Constructing a *Timestamp* or *DatetimeIndex* with an epoch timestamp with the `tz` argument specified will currently localize the epoch timestamps to UTC first then convert the result to the specified time zone. However, this behavior is *deprecated*, and if you have epochs in wall time in another timezone, it is recommended to read the epochs as timezone-naive timestamps and then localize to the appropriate timezone:

```
In [60]: pd.Timestamp(1262347200000000000).tz_localize('US/Pacific')
Out[60]: Timestamp('2010-01-01 12:00:00-0800', tz='US/Pacific')

In [61]: pd.DatetimeIndex([1262347200000000000]).tz_localize('US/Pacific')
Out[61]: DatetimeIndex(['2010-01-01 12:00:00-08:00'], dtype='datetime64[ns, US/
→Pacific]', freq=None)
```

---

**Note:** Epoch times will be rounded to the nearest nanosecond.

---

**Warning:** Conversion of float epoch times can lead to inaccurate and unexpected results. Python floats have about 15 digits precision in decimal. Rounding during conversion from float to high precision `Timestamp` is unavoidable. The only way to achieve exact precision is to use a fixed-width types (e.g. an int64).

```
In [62]: pd.to_datetime([1490195805.433, 1490195805.433502912], unit='s')
Out[62]: DatetimeIndex(['2017-03-22 15:16:45.433000088', '2017-03-22 15:16:45.
→433502913'], dtype='datetime64[ns]', freq=None)

In [63]: pd.to_datetime(1490195805433502912, unit='ns')
Out[63]: Timestamp('2017-03-22 15:16:45.433502912')
```

See also:

*Using the origin Parameter*

**From timestamps to epoch**

To invert the operation from above, namely, to convert from a `Timestamp` to a 'unix' epoch:

```
In [64]: stamps = pd.date_range('2012-10-08 18:15:05', periods=4, freq='D')

In [65]: stamps
Out[65]:
DatetimeIndex(['2012-10-08 18:15:05', '2012-10-09 18:15:05',
               '2012-10-10 18:15:05', '2012-10-11 18:15:05'],
              dtype='datetime64[ns]', freq='D')
```

We subtract the epoch (midnight at January 1, 1970 UTC) and then floor divide by the "unit" (1 second).

```
In [66]: (stamps - pd.Timestamp("1970-01-01")) // pd.Timedelta('1s')
Out[66]: Int64Index([1349720105, 1349806505, 1349892905, 1349979305], dtype='int64')
```

**Using the `origin` Parameter**

Using the `origin` parameter, one can specify an alternative starting point for creation of a `DatetimeIndex`. For example, to use 1960-01-01 as the starting date:

```
In [67]: pd.to_datetime([1, 2, 3], unit='D', origin=pd.Timestamp('1960-01-01'))
Out[67]: DatetimeIndex(['1960-01-02', '1960-01-03', '1960-01-04'], dtype=
→'datetime64[ns]', freq=None)
```

The default is set at `origin='unix'`, which defaults to `1970-01-01 00:00:00`. Commonly called 'unix epoch' or POSIX time.

```
In [68]: pd.to_datetime([1, 2, 3], unit='D')
Out[68]: DatetimeIndex(['1970-01-02', '1970-01-03', '1970-01-04'], dtype=
→'datetime64[ns]', freq=None)
```

## 2.14.4 Generating ranges of timestamps

To generate an index with timestamps, you can use either the `DatetimeIndex` or `Index` constructor and pass in a list of datetime objects:

```
In [69]: dates = [datetime.datetime(2012, 5, 1),
   ....:          datetime.datetime(2012, 5, 2),
   ....:          datetime.datetime(2012, 5, 3)]
   ....:

# Note the frequency information
In [70]: index = pd.DatetimeIndex(dates)

In [71]: index
Out[71]: DatetimeIndex(['2012-05-01', '2012-05-02', '2012-05-03'], dtype=
→'datetime64[ns]', freq=None)

# Automatically converted to DatetimeIndex
In [72]: index = pd.Index(dates)
```

(continues on next page)

```
In [73]: index
Out[73]: DatetimeIndex(['2012-05-01', '2012-05-02', '2012-05-03'], dtype=
→'datetime64[ns]', freq=None)
```

In practice 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 *date_range()* and *bdate_range()* functions to create a DatetimeIndex. The default frequency for date_range is a **calendar day** while the default for bdate_range is a **business day**:

```
In [74]: start = datetime.datetime(2011, 1, 1)

In [75]: end = datetime.datetime(2012, 1, 1)

In [76]: index = pd.date_range(start, end)

In [77]: index
Out[77]:
DatetimeIndex(['2011-01-01', '2011-01-02', '2011-01-03', '2011-01-04',
               '2011-01-05', '2011-01-06', '2011-01-07', '2011-01-08',
               '2011-01-09', '2011-01-10',
               ...
               '2011-12-23', '2011-12-24', '2011-12-25', '2011-12-26',
               '2011-12-27', '2011-12-28', '2011-12-29', '2011-12-30',
               '2011-12-31', '2012-01-01'],
              dtype='datetime64[ns]', length=366, freq='D')

In [78]: index = pd.bdate_range(start, end)

In [79]: index
Out[79]:
DatetimeIndex(['2011-01-03', '2011-01-04', '2011-01-05', '2011-01-06',
               '2011-01-07', '2011-01-10', '2011-01-11', '2011-01-12',
               '2011-01-13', '2011-01-14',
               ...
               '2011-12-19', '2011-12-20', '2011-12-21', '2011-12-22',
               '2011-12-23', '2011-12-26', '2011-12-27', '2011-12-28',
               '2011-12-29', '2011-12-30'],
              dtype='datetime64[ns]', length=260, freq='B')
```

Convenience functions like date_range and bdate_range can utilize a variety of *frequency aliases*:

```
In [80]: pd.date_range(start, periods=1000, freq='M')
Out[80]:
DatetimeIndex(['2011-01-31', '2011-02-28', '2011-03-31', '2011-04-30',
               '2011-05-31', '2011-06-30', '2011-07-31', '2011-08-31',
               '2011-09-30', '2011-10-31',
               ...
               '2093-07-31', '2093-08-31', '2093-09-30', '2093-10-31',
               '2093-11-30', '2093-12-31', '2094-01-31', '2094-02-28',
               '2094-03-31', '2094-04-30'],
              dtype='datetime64[ns]', length=1000, freq='M')

In [81]: pd.bdate_range(start, periods=250, freq='BQS')
Out[81]:
DatetimeIndex(['2011-01-03', '2011-04-01', '2011-07-01', '2011-10-03',
               '2012-01-02', '2012-04-02', '2012-07-02', '2012-10-01',
```

```
                '2013-01-01', '2013-04-01',
                ...
                '2071-01-01', '2071-04-01', '2071-07-01', '2071-10-01',
                '2072-01-01', '2072-04-01', '2072-07-01', '2072-10-03',
                '2073-01-02', '2073-04-03'],
               dtype='datetime64[ns]', length=250, freq='BQS-JAN')
```

`date_range` and `bdate_range` make it easy to generate a range of dates using various combinations of parameters like `start`, `end`, `periods`, and `freq`. The start and end dates are strictly inclusive, so dates outside of those specified will not be generated:

```
In [82]: pd.date_range(start, end, freq='BM')
Out[82]:
DatetimeIndex(['2011-01-31', '2011-02-28', '2011-03-31', '2011-04-29',
               '2011-05-31', '2011-06-30', '2011-07-29', '2011-08-31',
               '2011-09-30', '2011-10-31', '2011-11-30', '2011-12-30'],
              dtype='datetime64[ns]', freq='BM')

In [83]: pd.date_range(start, end, freq='W')
Out[83]:
DatetimeIndex(['2011-01-02', '2011-01-09', '2011-01-16', '2011-01-23',
               '2011-01-30', '2011-02-06', '2011-02-13', '2011-02-20',
               '2011-02-27', '2011-03-06', '2011-03-13', '2011-03-20',
               '2011-03-27', '2011-04-03', '2011-04-10', '2011-04-17',
               '2011-04-24', '2011-05-01', '2011-05-08', '2011-05-15',
               '2011-05-22', '2011-05-29', '2011-06-05', '2011-06-12',
               '2011-06-19', '2011-06-26', '2011-07-03', '2011-07-10',
               '2011-07-17', '2011-07-24', '2011-07-31', '2011-08-07',
               '2011-08-14', '2011-08-21', '2011-08-28', '2011-09-04',
               '2011-09-11', '2011-09-18', '2011-09-25', '2011-10-02',
               '2011-10-09', '2011-10-16', '2011-10-23', '2011-10-30',
               '2011-11-06', '2011-11-13', '2011-11-20', '2011-11-27',
               '2011-12-04', '2011-12-11', '2011-12-18', '2011-12-25',
               '2012-01-01'],
              dtype='datetime64[ns]', freq='W-SUN')

In [84]: pd.bdate_range(end=end, periods=20)
Out[84]:
DatetimeIndex(['2011-12-05', '2011-12-06', '2011-12-07', '2011-12-08',
               '2011-12-09', '2011-12-12', '2011-12-13', '2011-12-14',
               '2011-12-15', '2011-12-16', '2011-12-19', '2011-12-20',
               '2011-12-21', '2011-12-22', '2011-12-23', '2011-12-26',
               '2011-12-27', '2011-12-28', '2011-12-29', '2011-12-30'],
              dtype='datetime64[ns]', freq='B')

In [85]: pd.bdate_range(start=start, periods=20)
Out[85]:
DatetimeIndex(['2011-01-03', '2011-01-04', '2011-01-05', '2011-01-06',
               '2011-01-07', '2011-01-10', '2011-01-11', '2011-01-12',
               '2011-01-13', '2011-01-14', '2011-01-17', '2011-01-18',
               '2011-01-19', '2011-01-20', '2011-01-21', '2011-01-24',
               '2011-01-25', '2011-01-26', '2011-01-27', '2011-01-28'],
              dtype='datetime64[ns]', freq='B')
```

New in version 0.23.0.

Specifying `start`, `end`, and `periods` will generate a range of evenly spaced dates from `start` to `end` inclusively,

with `periods` number of elements in the resulting `DatetimeIndex`:

```
In [86]: pd.date_range('2018-01-01', '2018-01-05', periods=5)
Out[86]:
DatetimeIndex(['2018-01-01', '2018-01-02', '2018-01-03', '2018-01-04',
               '2018-01-05'],
              dtype='datetime64[ns]', freq=None)

In [87]: pd.date_range('2018-01-01', '2018-01-05', periods=10)
Out[87]:
DatetimeIndex(['2018-01-01 00:00:00', '2018-01-01 10:40:00',
               '2018-01-01 21:20:00', '2018-01-02 08:00:00',
               '2018-01-02 18:40:00', '2018-01-03 05:20:00',
               '2018-01-03 16:00:00', '2018-01-04 02:40:00',
               '2018-01-04 13:20:00', '2018-01-05 00:00:00'],
              dtype='datetime64[ns]', freq=None)
```

### Custom frequency ranges

`bdate_range` can also generate a range of custom frequency dates by using the `weekmask` and `holidays` parameters. These parameters will only be used if a custom frequency string is passed.

```
In [88]: weekmask = 'Mon Wed Fri'

In [89]: holidays = [datetime.datetime(2011, 1, 5), datetime.datetime(2011, 3, 14)]

In [90]: pd.bdate_range(start, end, freq='C', weekmask=weekmask, holidays=holidays)
Out[90]:
DatetimeIndex(['2011-01-03', '2011-01-07', '2011-01-10', '2011-01-12',
               '2011-01-14', '2011-01-17', '2011-01-19', '2011-01-21',
               '2011-01-24', '2011-01-26',
               ...
               '2011-12-09', '2011-12-12', '2011-12-14', '2011-12-16',
               '2011-12-19', '2011-12-21', '2011-12-23', '2011-12-26',
               '2011-12-28', '2011-12-30'],
              dtype='datetime64[ns]', length=154, freq='C')

In [91]: pd.bdate_range(start, end, freq='CBMS', weekmask=weekmask)
Out[91]:
DatetimeIndex(['2011-01-03', '2011-02-02', '2011-03-02', '2011-04-01',
               '2011-05-02', '2011-06-01', '2011-07-01', '2011-08-01',
               '2011-09-02', '2011-10-03', '2011-11-02', '2011-12-02'],
              dtype='datetime64[ns]', freq='CBMS')
```

See also:

*Custom business days*

### 2.14.5 Timestamp limitations

Since pandas represents timestamps in nanosecond resolution, the time span that can be represented using a 64-bit integer is limited to approximately 584 years:

```
In [92]: pd.Timestamp.min
Out[92]: Timestamp('1677-09-21 00:12:43.145225')

In [93]: pd.Timestamp.max
Out[93]: Timestamp('2262-04-11 23:47:16.854775807')
```

**See also:**

*Representing out-of-bounds spans*

### 2.14.6 Indexing

One of the main uses for `DatetimeIndex` is as an index for pandas objects. The `DatetimeIndex` class contains many time series 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` objects have all the basic functionality of regular `Index` objects, and a smorgasbord of advanced time series 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.

---

`DatetimeIndex` can be used like a regular index and offers all of its intelligent functionality like selection, slicing, etc.

```
In [94]: rng = pd.date_range(start, end, freq='BM')

In [95]: ts = pd.Series(np.random.randn(len(rng)), index=rng)

In [96]: ts.index
Out[96]:
DatetimeIndex(['2011-01-31', '2011-02-28', '2011-03-31', '2011-04-29',
               '2011-05-31', '2011-06-30', '2011-07-29', '2011-08-31',
               '2011-09-30', '2011-10-31', '2011-11-30', '2011-12-30'],
              dtype='datetime64[ns]', freq='BM')

In [97]: ts[:5].index
Out[97]:
```

(continues on next page)

```
DatetimeIndex(['2011-01-31', '2011-02-28', '2011-03-31', '2011-04-29',
               '2011-05-31'],
              dtype='datetime64[ns]', freq='BM')

In [98]: ts[::2].index
Out[98]:
DatetimeIndex(['2011-01-31', '2011-03-31', '2011-05-31', '2011-07-29',
               '2011-09-30', '2011-11-30'],
              dtype='datetime64[ns]', freq='2BM')
```

### Partial string indexing

Dates and strings that parse to timestamps can be passed as indexing parameters:

```
In [99]: ts['1/31/2011']
Out[99]: 0.11920871129693428

In [100]: ts[datetime.datetime(2011, 12, 25):]
Out[100]:
2011-12-30    0.56702
Freq: BM, dtype: float64

In [101]: ts['10/31/2011':'12/31/2011']
Out[101]:
2011-10-31    0.271860
2011-11-30   -0.424972
2011-12-30    0.567020
Freq: BM, dtype: float64
```

To provide convenience for accessing longer time series, you can also pass in the year or year and month as strings:

```
In [102]: ts['2011']
Out[102]:
2011-01-31    0.119209
2011-02-28   -1.044236
2011-03-31   -0.861849
2011-04-29   -2.104569
2011-05-31   -0.494929
2011-06-30    1.071804
2011-07-29    0.721555
2011-08-31   -0.706771
2011-09-30   -1.039575
2011-10-31    0.271860
2011-11-30   -0.424972
2011-12-30    0.567020
Freq: BM, dtype: float64

In [103]: ts['2011-6']
Out[103]:
2011-06-30    1.071804
Freq: BM, dtype: float64
```

This type of slicing will work on a `DataFrame` with a `DatetimeIndex` as well. Since the partial string selection is a form of label slicing, the endpoints **will be** included. This would include matching times on an included date:

```
In [104]: dft = pd.DataFrame(np.random.randn(100000, 1), columns=['A'],
   .....:                    index=pd.date_range('20130101', periods=100000, freq='T
→'))
   .....:

In [105]: dft
Out[105]:
                            A
2013-01-01 00:00:00  0.276232
2013-01-01 00:01:00 -1.087401
2013-01-01 00:02:00 -0.673690
2013-01-01 00:03:00  0.113648
2013-01-01 00:04:00 -1.478427
...                       ...
2013-03-11 10:35:00 -0.747967
2013-03-11 10:36:00 -0.034523
2013-03-11 10:37:00 -0.201754
2013-03-11 10:38:00 -1.509067
2013-03-11 10:39:00 -1.693043

[100000 rows x 1 columns]

In [106]: dft['2013']
Out[106]:
                            A
2013-01-01 00:00:00  0.276232
2013-01-01 00:01:00 -1.087401
2013-01-01 00:02:00 -0.673690
2013-01-01 00:03:00  0.113648
2013-01-01 00:04:00 -1.478427
...                       ...
2013-03-11 10:35:00 -0.747967
2013-03-11 10:36:00 -0.034523
2013-03-11 10:37:00 -0.201754
2013-03-11 10:38:00 -1.509067
2013-03-11 10:39:00 -1.693043

[100000 rows x 1 columns]
```

This starts on the very first time in the month, and includes the last date and time for the month:

```
In [107]: dft['2013-1':'2013-2']
Out[107]:
                            A
2013-01-01 00:00:00  0.276232
2013-01-01 00:01:00 -1.087401
2013-01-01 00:02:00 -0.673690
2013-01-01 00:03:00  0.113648
2013-01-01 00:04:00 -1.478427
...                       ...
2013-02-28 23:55:00  0.850929
2013-02-28 23:56:00  0.976712
2013-02-28 23:57:00 -2.693884
2013-02-28 23:58:00 -1.575535
2013-02-28 23:59:00 -1.573517

[84960 rows x 1 columns]
```