

- azuranski +
- behzad nouri
- cel4
- emilydolson +
- hironow +
- lexical
- llllllllll +
- rockg
- silentquasar +
- sinhrks
- taeold +

5.10.2 v0.17.0 (October 9, 2015)

This is a major release from 0.16.2 and includes a small number of API changes, several new features, enhancements, and performance improvements along with a large number of bug fixes. We recommend that all users upgrade to this version.

Warning: pandas \geq 0.17.0 will no longer support compatibility with Python version 3.2 ([GH9118](#))

Warning: The `pandas.io.data` package is deprecated and will be replaced by the `pandas-datareader` package. This will allow the data modules to be independently updated to your pandas installation. The API for `pandas-datareader v0.1.1` is exactly the same as in pandas v0.17.0 ([GH8961](#), [GH10861](#)).

After installing `pandas-datareader`, you can easily change your imports:

```
from pandas.io import data, wb
```

becomes

```
from pandas_datareader import data, wb
```

Highlights include:

- Release the Global Interpreter Lock (GIL) on some cython operations, see [here](#)
- Plotting methods are now available as attributes of the `.plot` accessor, see [here](#)
- The sorting API has been revamped to remove some long-time inconsistencies, see [here](#)
- Support for a `datetime64[ns]` with timezones as a first-class dtype, see [here](#)
- The default for `to_datetime` will now be to `raise` when presented with unparseable formats, previously this would return the original input. Also, date parse functions now return consistent results. See [here](#)
- The default for `dropna` in `HDFStore` has changed to `False`, to store by default all rows even if they are all NaN, see [here](#)
- Datetime accessor (`dt`) now supports `Series.dt.strftime` to generate formatted strings for datetime-likes, and `Series.dt.total_seconds` to generate each duration of the timedelta in seconds. See [here](#)

- `Period` and `PeriodIndex` can handle multiplied freq like 3D, which corresponding to 3 days span. See [here](#)
- Development installed versions of pandas will now have PEP440 compliant version strings ([GH9518](#))
- Development support for benchmarking with the [Air Speed Velocity](#) library ([GH8361](#))
- Support for reading SAS xport files, see [here](#)
- Documentation comparing SAS to *pandas*, see [here](#)
- Removal of the automatic `TimeSeries` broadcasting, deprecated since 0.8.0, see [here](#)
- Display format with plain text can optionally align with Unicode East Asian Width, see [here](#)
- Compatibility with Python 3.5 ([GH11097](#))
- Compatibility with matplotlib 1.5.0 ([GH11111](#))

Check the [API Changes](#) and [deprecations](#) before updating.

What's new in v0.17.0

- *New features*
 - *Datetime with TZ*
 - *Releasing the GIL*
 - *Plot submethods*
 - *Additional methods for dt accessor*
 - * *strftime*
 - * *total_seconds*
 - *Period frequency enhancement*
 - *Support for SAS XPORT files*
 - *Support for math functions in .eval()*
 - *Changes to Excel with MultiIndex*
 - *Google BigQuery enhancements*
 - *Display alignment with Unicode East Asian width*
 - *Other enhancements*
- *Backwards incompatible API changes*
 - *Changes to sorting API*
 - *Changes to to_datetime and to_timedelta*
 - * *Error handling*
 - * *Consistent parsing*
 - *Changes to Index comparisons*
 - *Changes to boolean comparisons vs. None*
 - *HDFStore dropna behavior*
 - *Changes to display.precision option*
 - *Changes to Categorical.unique*

- *Changes to `bool` passed as header in parsers*
- *Other API changes*
- *Deprecations*
- *Removal of prior version deprecations/changes*
- *Performance improvements*
- *Bug fixes*
- *Contributors*

New features

Datetime with TZ

We are adding an implementation that natively supports datetime with timezones. A `Series` or a `DataFrame` column previously *could* be assigned a datetime with timezones, and would work as an object dtype. This had performance issues with a large number rows. See the [docs](#) for more details. ([GH8260](#), [GH10763](#), [GH11034](#)).

The new implementation allows for having a single-timezone across all rows, with operations in a performant manner.

```
In [1]: df = pd.DataFrame({'A': pd.date_range('20130101', periods=3),
...:                      'B': pd.date_range('20130101', periods=3, tz='US/Eastern'),
...:                      'C': pd.date_range('20130101', periods=3, tz='CET')})
...:
```

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

```
Out[2]:
```

	A	B	C
0	2013-01-01 2013-01-01 00:00:00-05:00	2013-01-01 00:00:00+01:00	
1	2013-01-02 2013-01-02 00:00:00-05:00	2013-01-02 00:00:00+01:00	
2	2013-01-03 2013-01-03 00:00:00-05:00	2013-01-03 00:00:00+01:00	

```
[3 rows x 3 columns]
```

```
In [3]: df.dtypes
```

```
Out[3]:
```

```
A          datetime64[ns]
B    datetime64[ns, US/Eastern]
C    datetime64[ns, CET]
Length: 3, dtype: object
```

```
In [4]: df.B
```

```
Out[4]:
```

```
0    2013-01-01 00:00:00-05:00
1    2013-01-02 00:00:00-05:00
2    2013-01-03 00:00:00-05:00
Name: B, Length: 3, dtype: datetime64[ns, US/Eastern]
```

```
In [5]: df.B.dt.tz_localize(None)
```

```
Out[5]:
```

```
0    2013-01-01
1    2013-01-02
2    2013-01-03
Name: B, Length: 3, dtype: datetime64[ns]
```

This uses a new-dtype representation as well, that is very similar in look-and-feel to its numpy cousin `datetime64[ns]`

```
In [6]: df['B'].dtype
Out[6]: datetime64[ns, US/Eastern]

In [7]: type(df['B'].dtype)
Out[7]: pandas.core.dtypes.dtypes.DatetimeTZDtype
```

Note: There is a slightly different string repr for the underlying `DatetimeIndex` as a result of the dtype changes, but functionally these are the same.

Previous behavior:

```
In [1]: pd.date_range('20130101', periods=3, tz='US/Eastern')
Out[1]: DatetimeIndex(['2013-01-01 00:00:00-05:00', '2013-01-02 00:00:00-05:00',
                        '2013-01-03 00:00:00-05:00'],
                        dtype='datetime64[ns]', freq='D', tz='US/Eastern')

In [2]: pd.date_range('20130101', periods=3, tz='US/Eastern').dtype
Out[2]: dtype('<M8[ns]')
```

New behavior:

```
In [8]: pd.date_range('20130101', periods=3, tz='US/Eastern')
Out[8]:
DatetimeIndex(['2013-01-01 00:00:00-05:00', '2013-01-02 00:00:00-05:00',
                '2013-01-03 00:00:00-05:00'],
                dtype='datetime64[ns, US/Eastern]', freq='D')

In [9]: pd.date_range('20130101', periods=3, tz='US/Eastern').dtype
Out[9]: datetime64[ns, US/Eastern]
```

Releasing the GIL

We are releasing the global-interpreter-lock (GIL) on some cython operations. This will allow other threads to run simultaneously during computation, potentially allowing performance improvements from multi-threading. Notably `groupby`, `nsmallest`, `value_counts` and some indexing operations benefit from this. ([GH8882](#))

For example the `groupby` expression in the following code will have the GIL released during the factorization step, e.g. `df.groupby('key')` as well as the `.sum()` operation.

```
N = 1000000
ngroups = 10
df = DataFrame({'key': np.random.randint(0, ngroups, size=N),
                'data': np.random.randn(N)})
df.groupby('key')['data'].sum()
```

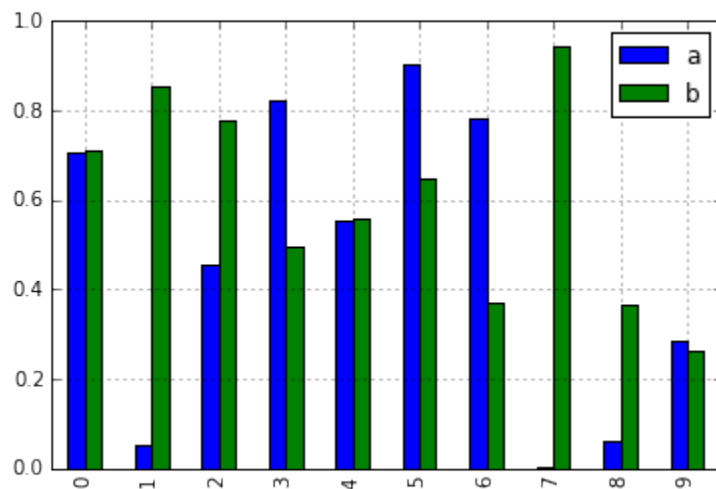
Releasing of the GIL could benefit an application that uses threads for user interactions (e.g. [QT](#)), or performing multi-threaded computations. A nice example of a library that can handle these types of computation-in-parallel is the [dask](#) library.

Plot submethods

The Series and DataFrame `.plot()` method allows for customizing *plot types* by supplying the `kind` keyword arguments. Unfortunately, many of these kinds of plots use different required and optional keyword arguments, which makes it difficult to discover what any given plot kind uses out of the dozens of possible arguments.

To alleviate this issue, we have added a new, optional plotting interface, which exposes each kind of plot as a method of the `.plot` attribute. Instead of writing `series.plot(kind=<kind>, ...)`, you can now also use `series.plot.<kind>(...)`:

```
In [10]: df = pd.DataFrame(np.random.rand(10, 2), columns=['a', 'b'])
In [11]: df.plot.bar()
```



As a result of this change, these methods are now all discoverable via tab-completion:

```
In [12]: df.plot.<TAB> # noqa: E225, E999
df.plot.area      df.plot.barh      df.plot.density  df.plot.hist      df.plot.line
↳ df.plot.scatter
df.plot.bar       df.plot.box       df.plot.hexbin   df.plot.kde       df.plot.pie
```

Each method signature only includes relevant arguments. Currently, these are limited to required arguments, but in the future these will include optional arguments, as well. For an overview, see the new *Plotting* API documentation.

Additional methods for `dt` accessor

strftime

We are now supporting a `Series.dt.strftime` method for datetime-likes to generate a formatted string (GH10110). Examples:

```
# DatetimeIndex
In [13]: s = pd.Series(pd.date_range('20130101', periods=4))

In [14]: s
Out[14]:
0    2013-01-01
```

(continues on next page)

(continued from previous page)

```

1    2013-01-02
2    2013-01-03
3    2013-01-04
Length: 4, dtype: datetime64[ns]

In [15]: s.dt.strftime('%Y/%m/%d')
Out[15]:
0    2013/01/01
1    2013/01/02
2    2013/01/03
3    2013/01/04
Length: 4, dtype: object

```

```

# PeriodIndex
In [16]: s = pd.Series(pd.period_range('20130101', periods=4))

In [17]: s
Out[17]:
0    2013-01-01
1    2013-01-02
2    2013-01-03
3    2013-01-04
Length: 4, dtype: period[D]

In [18]: s.dt.strftime('%Y/%m/%d')
Out[18]:
0    2013/01/01
1    2013/01/02
2    2013/01/03
3    2013/01/04
Length: 4, dtype: object

```

The string format is as the python standard library and details can be found [here](#)

total_seconds

`pd.Series` of type `timedelta64` has new method `.dt.total_seconds()` returning the duration of the `timedelta` in seconds ([GH10817](#))

```

# TimedeltaIndex
In [19]: s = pd.Series(pd.timedelta_range('1 minutes', periods=4))

In [20]: s
Out[20]:
0    0 days 00:01:00
1    1 days 00:01:00
2    2 days 00:01:00
3    3 days 00:01:00
Length: 4, dtype: timedelta64[ns]

In [21]: s.dt.total_seconds()
Out[21]:
0         60.0
1      86460.0
2     172860.0

```

(continues on next page)

(continued from previous page)

```
3      259260.0
Length: 4, dtype: float64
```

Period frequency enhancement

`Period`, `PeriodIndex` and `period_range` can now accept multiplied freq. Also, `Period.freq` and `PeriodIndex.freq` are now stored as a `DateOffset` instance like `DatetimeIndex`, and not as `str` (GH7811)

A multiplied freq represents a span of corresponding length. The example below creates a period of 3 days. Addition and subtraction will shift the period by its span.

```
In [22]: p = pd.Period('2015-08-01', freq='3D')

In [23]: p
Out[23]: Period('2015-08-01', '3D')

In [24]: p + 1
Out[24]: Period('2015-08-04', '3D')

In [25]: p - 2
Out[25]: Period('2015-07-26', '3D')

In [26]: p.to_timestamp()
Out[26]: Timestamp('2015-08-01 00:00:00')

In [27]: p.to_timestamp(how='E')
Out[27]: Timestamp('2015-08-03 23:59:59.999999999')
```

You can use the multiplied freq in `PeriodIndex` and `period_range`.

```
In [28]: idx = pd.period_range('2015-08-01', periods=4, freq='2D')

In [29]: idx
Out[29]: PeriodIndex(['2015-08-01', '2015-08-03', '2015-08-05', '2015-08-07'], dtype=
↪ 'period[2D]', freq='2D')

In [30]: idx + 1
Out[30]: PeriodIndex(['2015-08-03', '2015-08-05', '2015-08-07', '2015-08-09'], dtype=
↪ 'period[2D]', freq='2D')
```

Support for SAS XPORT files

`read_sas()` provides support for reading *SAS XPORT* format files. (GH4052).

```
df = pd.read_sas('sas_xport.xpt')
```

It is also possible to obtain an iterator and read an XPORT file incrementally.

```
for df in pd.read_sas('sas_xport.xpt', chunksize=10000):
    do_something(df)
```

See the [docs](#) for more details.

Support for math functions in .eval()

`eval()` now supports calling math functions ([GH4893](#))

```
df = pd.DataFrame({'a': np.random.randn(10)})
df.eval("b = sin(a)")
```

The support math functions are *sin*, *cos*, *exp*, *log*, *expm1*, *log1p*, *sqrt*, *sinh*, *cosh*, *tanh*, *arcsin*, *arccos*, *arctan*, *arccosh*, *arcsinh*, *arctanh*, *abs* and *arctan2*.

These functions map to the intrinsics for the NumExpr engine. For the Python engine, they are mapped to NumPy calls.

Changes to Excel with MultiIndex

In version 0.16.2 a `DataFrame` with `MultiIndex` columns could not be written to Excel via `to_excel`. That functionality has been added ([GH10564](#)), along with updating `read_excel` so that the data can be read back with, no loss of information, by specifying which columns/rows make up the `MultiIndex` in the header and `index_col` parameters ([GH4679](#))

See the [documentation](#) for more details.

```
In [31]: df = pd.DataFrame([[1, 2, 3, 4], [5, 6, 7, 8]],
.....:                     columns=pd.MultiIndex.from_product(
.....:                     [['foo', 'bar'], ['a', 'b']], names=['col1', 'col2']),
.....:                     index=pd.MultiIndex.from_product(['j', 'l', 'k'],
.....:                                                         names=['i1', 'i2']))
.....:

In [32]: df
Out[32]:
col1 foo    bar
col2  a  b   a  b
i1 i2
j  l    1  2   3  4
   k    5  6   7  8

[2 rows x 4 columns]

In [33]: df.to_excel('test.xlsx')

In [34]: df = pd.read_excel('test.xlsx', header=[0, 1], index_col=[0, 1])

In [35]: df
Out[35]:
col1 foo    bar
col2  a  b   a  b
i1 i2
j  l    1  2   3  4
   k    5  6   7  8

[2 rows x 4 columns]
```

Previously, it was necessary to specify the `has_index_names` argument in `read_excel`, if the serialized data had index names. For version 0.17.0 the output format of `to_excel` has been changed to make this keyword unnecessary - the change is shown below.

Old

	A	B	C	D	E	F
1		A	B	C	D	
2	idx_name					
3	2000-01-07 00:00:00	0.968129	0.906529	0.05343	0.02619	
4	2000-01-10 00:00:00	-0.16632	1.981993	1.833093	0.803685	
5	2000-01-11 00:00:00	0.121057	0.36946	-0.02888	1.683975	
6	2000-01-12 00:00:00	-1.70456	-0.73098	-0.38088	0.020946	
7	2000-01-13 00:00:00	-1.20024	1.907733	0.629318	1.507033	
8	2000-01-14 00:00:00	-0.66344	0.073188	1.583482	0.735205	
9	2000-01-17 00:00:00	0.716635	-2.07952	1.760536	0.970309	

New

	A	B	C	D	E	
1	idx_name	A	B	C	D	
2	2000-01-07 00:00:00	0.968129	0.906529	0.05343	0.02619	
3	2000-01-10 00:00:00	-0.16632	1.981993	1.833093	0.803685	
4	2000-01-11 00:00:00	0.121057	0.36946	-0.02888	1.683975	
5	2000-01-12 00:00:00	-1.70456	-0.73098	-0.38088	0.020946	
6	2000-01-13 00:00:00	-1.20024	1.907733	0.629318	1.507033	
7	2000-01-14 00:00:00	-0.66344	0.073188	1.583482	0.735205	
8	2000-01-17 00:00:00	0.716635	-2.07952	1.760536	0.970309	
9	2000-01-18 00:00:00	0.727629	2.22267	2.706276	0.681942	

Warning: Excel files saved in version 0.16.2 or prior that had index names will still be able to be read in, but the `has_index_names` argument must be specified to `True`.

Google BigQuery enhancements

- Added ability to automatically create a table/dataset using the `pandas.io.gbq.to_gbq()` function if the destination table/dataset does not exist. (GH8325, GH11121).
- Added ability to replace an existing table and schema when calling the `pandas.io.gbq.to_gbq()` function via the `if_exists` argument. See the docs for more details (GH8325).
- `InvalidColumnOrder` and `InvalidPageToken` in the `gbq` module will raise `ValueError` instead of `IOError`.
- The `generate_bq_schema()` function is now deprecated and will be removed in a future version (GH11121).
- The `gbq` module will now support Python 3 (GH11094).

Display alignment with Unicode East Asian width

Warning: Enabling this option will affect the performance for printing of `DataFrame` and `Series` (about 2 times slower). Use only when it is actually required.

Some East Asian countries use Unicode characters its width is corresponding to 2 alphabets. If a `DataFrame` or `Series` contains these characters, the default output cannot be aligned properly. The following options are added to enable precise handling for these characters.

- `display.unicode.east_asian_width`: Whether to use the Unicode East Asian Width to calculate the display text width. (GH2612)
- `display.unicode.ambiguous_as_wide`: Whether to handle Unicode characters belong to Ambiguous as Wide. (GH11102)

```
In [36]: df = pd.DataFrame({'u': ['UK', u''], u': ['Alice', u'']})
```

```
In [37]: df;
```

```
>>> df = pd.DataFrame({'u'国籍': ['UK', u'日本'], u'名前': ['Alice', u'しのぶ']})
>>> df
      名前  国籍
0  Alice  UK
1   し の ぶ  日本
```

```
In [38]: pd.set_option('display.unicode.east_asian_width', True)
```

```
In [39]: df;
```

```
>>> pd.set_option('display.unicode.east_asian_width', True)
>>> df
      名前  国籍
0  Alice  UK
1   し の ぶ  日本
```

For further details, see [here](#)

Other enhancements

- Support for `openpyxl >= 2.2`. The API for style support is now stable (GH10125)
- `merge` now accepts the argument `indicator` which adds a Categorical-type column (by default called `_merge`) to the output object that takes on the values (GH8790)

Observation Origin	<code>_merge</code> value
Merge key only in 'left' frame	<code>left_only</code>
Merge key only in 'right' frame	<code>right_only</code>
Merge key in both frames	<code>both</code>

```
In [40]: df1 = pd.DataFrame({'col1': [0, 1], 'col_left': ['a', 'b']})
```

```
In [41]: df2 = pd.DataFrame({'col1': [1, 2, 2], 'col_right': [2, 2, 2]})
```

(continues on next page)

(continued from previous page)

```
In [42]: pd.merge(df1, df2, on='coll', how='outer', indicator=True)
Out[42]:
```

	coll	col_left	col_right	_merge
0	0	a	NaN	left_only
1	1	b	2.0	both
2	2	NaN	2.0	right_only
3	2	NaN	2.0	right_only

```
[4 rows x 4 columns]
```

For more, see the [updated docs](#)

- `pd.to_numeric` is a new function to coerce strings to numbers (possibly with coercion) ([GH11133](#))
- `pd.merge` will now allow duplicate column names if they are not merged upon ([GH10639](#)).
- `pd.pivot` will now allow passing index as None ([GH3962](#)).
- `pd.concat` will now use existing Series names if provided ([GH10698](#)).

```
In [43]: foo = pd.Series([1, 2], name='foo')
In [44]: bar = pd.Series([1, 2])
In [45]: baz = pd.Series([4, 5])
```

Previous behavior:

```
In [1]: pd.concat([foo, bar, baz], 1)
Out[1]:
```

	0	1	2
0	1	1	4
1	2	2	5

New behavior:

```
In [46]: pd.concat([foo, bar, baz], 1)
Out[46]:
```

	foo	0	1
0	1	1	4
1	2	2	5

```
[2 rows x 3 columns]
```

- `DataFrame` has gained the `nlargest` and `nsmallest` methods ([GH10393](#))
- Add a `limit_direction` keyword argument that works with `limit` to enable interpolate to fill NaN values forward, backward, or both ([GH9218](#), [GH10420](#), [GH11115](#))

```
In [47]: ser = pd.Series([np.nan, np.nan, 5, np.nan, np.nan, np.nan, 13])
In [48]: ser.interpolate(limit=1, limit_direction='both')
Out[48]:
```

0	NaN
1	5.0
2	5.0
3	7.0
4	NaN

(continues on next page)

(continued from previous page)

```

5      11.0
6      13.0
Length: 7, dtype: float64

```

- Added a `DataFrame.round` method to round the values to a variable number of decimal places ([GH10568](#)).

```

In [49]: df = pd.DataFrame(np.random.random([3, 3]),
.....:                      columns=['A', 'B', 'C'],
.....:                      index=['first', 'second', 'third'])
.....:

```

```
In [50]: df
```

```

Out [50]:
           A          B          C
first  0.126970  0.966718  0.260476
second 0.897237  0.376750  0.336222
third   0.451376  0.840255  0.123102

```

```
[3 rows x 3 columns]
```

```
In [51]: df.round(2)
```

```

Out [51]:
           A          B          C
first   0.13   0.97   0.26
second  0.90   0.38   0.34
third   0.45   0.84   0.12

```

```
[3 rows x 3 columns]
```

```
In [52]: df.round({'A': 0, 'C': 2})
```

```

Out [52]:
           A          B          C
first   0.0   0.966718   0.26
second  1.0   0.376750   0.34
third   0.0   0.840255   0.12

```

```
[3 rows x 3 columns]
```

- `drop_duplicates` and `duplicated` now accept a `keep` keyword to target first, last, and all duplicates. The `take_last` keyword is deprecated, see [here](#) ([GH6511](#), [GH8505](#))

```
In [53]: s = pd.Series(['A', 'B', 'C', 'A', 'B', 'D'])
```

```
In [54]: s.drop_duplicates()
```

```

Out [54]:
0      A
1      B
2      C
5      D
Length: 4, dtype: object

```

```
In [55]: s.drop_duplicates(keep='last')
```

```

Out [55]:
2      C
3      A
4      B

```

(continues on next page)

(continued from previous page)

```

5      D
Length: 4, dtype: object

In [56]: s.drop_duplicates(keep=False)
Out[56]:
2      C
5      D
Length: 2, dtype: object

```

- Reindex now has a `tolerance` argument that allows for finer control of *Limits on filling while reindexing* (GH10411):

```

In [57]: df = pd.DataFrame({'x': range(5),
.....:                    't': pd.date_range('2000-01-01', periods=5)})
.....:

In [58]: df.reindex([0.1, 1.9, 3.5],
.....:              method='nearest',
.....:              tolerance=0.2)
.....:
Out[58]:
      x      t
0.1  0.0 2000-01-01
1.9  2.0 2000-01-03
3.5  NaN      NaT

[3 rows x 2 columns]

```

When used on a `DatetimeIndex`, `TimedeltaIndex` or `PeriodIndex`, `tolerance` will be coerced into a `Timedelta` if possible. This allows you to specify tolerance with a string:

```

In [59]: df = df.set_index('t')

In [60]: df.reindex(pd.to_datetime(['1999-12-31']),
.....:              method='nearest',
.....:              tolerance='1 day')
.....:
Out[60]:
      x
1999-12-31  0

[1 rows x 1 columns]

```

`tolerance` is also exposed by the lower level `Index.get_indexer` and `Index.get_loc` methods.

- Added functionality to use the `base` argument when resampling a `TimeDeltaIndex` (GH10530)
- `DatetimeIndex` can be instantiated using strings that contain `NaT` (GH7599)
- `to_datetime` can now accept the `yearfirst` keyword (GH7599)
- `pandas.tseries.offsets` larger than the `Day` offset can now be used with a `Series` for addition/subtraction (GH10699). See the *docs* for more details.
- `pd.Timedelta.total_seconds()` now returns `Timedelta` duration to ns precision (previously microsecond precision) (GH10939)
- `PeriodIndex` now supports arithmetic with `np.ndarray` (GH10638)

- Support pickling of `Period` objects ([GH10439](#))
- `.as_blocks` will now take a `copy` optional argument to return a copy of the data, default is to copy (no change in behavior from prior versions), ([GH9607](#))
- `regex` argument to `DataFrame.filter` now handles numeric column names instead of raising `ValueError` ([GH10384](#)).
- Enable reading gzip compressed files via URL, either by explicitly setting the compression parameter or by inferring from the presence of the HTTP Content-Encoding header in the response ([GH8685](#))
- Enable writing Excel files in *memory* using `StringIO/BytesIO` ([GH7074](#))
- Enable serialization of lists and dicts to strings in `ExcelWriter` ([GH8188](#))
- SQL io functions now accept a SQLAlchemy connectable. ([GH7877](#))
- `pd.read_sql` and `to_sql` can accept database URI as `con` parameter ([GH10214](#))
- `read_sql_table` will now allow reading from views ([GH10750](#)).
- Enable writing complex values to `HDFStores` when using the `table` format ([GH10447](#))
- Enable `pd.read_hdf` to be used without specifying a key when the HDF file contains a single dataset ([GH10443](#))
- `pd.read_stata` will now read Stata 118 type files. ([GH9882](#))
- `msgpack` submodule has been updated to 0.4.6 with backward compatibility ([GH10581](#))
- `DataFrame.to_dict` now accepts `orient='index'` keyword argument ([GH10844](#)).
- `DataFrame.apply` will return a `Series` of dicts if the passed function returns a dict and `reduce=True` ([GH8735](#)).
- Allow passing *kwargs* to the interpolation methods ([GH10378](#)).
- Improved error message when concatenating an empty iterable of `Dataframe` objects ([GH9157](#))
- `pd.read_csv` can now read bz2-compressed files incrementally, and the C parser can read bz2-compressed files from AWS S3 ([GH11070](#), [GH11072](#)).
- In `pd.read_csv`, recognize `s3n://` and `s3a://` URLs as designating S3 file storage ([GH11070](#), [GH11071](#)).
- Read CSV files from AWS S3 incrementally, instead of first downloading the entire file. (Full file download still required for compressed files in Python 2.) ([GH11070](#), [GH11073](#))
- `pd.read_csv` is now able to infer compression type for files read from AWS S3 storage ([GH11070](#), [GH11074](#)).

Backwards incompatible API changes

Changes to sorting API

The sorting API has had some longtime inconsistencies. ([GH9816](#), [GH8239](#)).

Here is a summary of the API **PRIOR** to 0.17.0:

- `Series.sort` is **INPLACE** while `DataFrame.sort` returns a new object.
- `Series.order` returns a new object
- It was possible to use `Series/DataFrame.sort_index` to sort by **values** by passing the `by` keyword.

- `Series/DataFrame.sortlevel` worked only on a `MultiIndex` for sorting by index.

To address these issues, we have revamped the API:

- We have introduced a new method, `DataFrame.sort_values()`, which is the merger of `DataFrame.sort()`, `Series.sort()`, and `Series.order()`, to handle sorting of **values**.
- The existing methods `Series.sort()`, `Series.order()`, and `DataFrame.sort()` have been deprecated and will be removed in a future version.
- The `by` argument of `DataFrame.sort_index()` has been deprecated and will be removed in a future version.
- The existing method `.sort_index()` will gain the `level` keyword to enable level sorting.

We now have two distinct and non-overlapping methods of sorting. A * marks items that will show a `FutureWarning`.

To sort by the **values**:

Previous	Replacement
* <code>Series.order()</code>	<code>Series.sort_values()</code>
* <code>Series.sort()</code>	<code>Series.sort_values(inplace=True)</code>
* <code>DataFrame.sort(columns=...)</code>	<code>DataFrame.sort_values(by=...)</code>

To sort by the **index**:

Previous	Replacement
<code>Series.sort_index()</code>	<code>Series.sort_index()</code>
<code>Series.sortlevel(level=...)</code>	<code>Series.sort_index(level=...)</code>
<code>DataFrame.sort_index()</code>	<code>DataFrame.sort_index()</code>
<code>DataFrame.sortlevel(level=...)</code>	<code>DataFrame.sort_index(level=...)</code>
* <code>DataFrame.sort()</code>	<code>DataFrame.sort_index()</code>

We have also deprecated and changed similar methods in two Series-like classes, `Index` and `Categorical`.

Previous	Replacement
* <code>Index.order()</code>	<code>Index.sort_values()</code>
* <code>Categorical.order()</code>	<code>Categorical.sort_values()</code>

Changes to `to_datetime` and `timedelta`

Error handling

The default for `pd.to_datetime` error handling has changed to `errors='raise'`. In prior versions it was `errors='ignore'`. Furthermore, the `coerce` argument has been deprecated in favor of `errors='coerce'`. This means that invalid parsing will raise rather than return the original input as in previous versions. ([GH10636](#))

Previous behavior:

```
In [2]: pd.to_datetime(['2009-07-31', 'asd'])
Out[2]: array(['2009-07-31', 'asd'], dtype=object)
```

New behavior:

```
In [3]: pd.to_datetime(['2009-07-31', 'asd'])
ValueError: Unknown string format
```

Of course you can coerce this as well.

```
In [61]: pd.to_datetime(['2009-07-31', 'asd'], errors='coerce')
Out [61]: DatetimeIndex(['2009-07-31', 'NaT'], dtype='datetime64[ns]', freq=None)
```

To keep the previous behavior, you can use `errors='ignore'`:

```
In [62]: pd.to_datetime(['2009-07-31', 'asd'], errors='ignore')
Out [62]: Index(['2009-07-31', 'asd'], dtype='object')
```

Furthermore, `pd.to_timedelta` has gained a similar API, of `errors='raise' | 'ignore' | 'coerce'`, and the `coerce` keyword has been deprecated in favor of `errors='coerce'`.

Consistent parsing

The string parsing of `to_datetime`, `Timestamp` and `DatetimeIndex` has been made consistent. ([GH7599](#))

Prior to v0.17.0, `Timestamp` and `to_datetime` may parse year-only datetime-string incorrectly using today's date, otherwise `DatetimeIndex` uses the beginning of the year. `Timestamp` and `to_datetime` may raise `ValueError` in some types of datetime-string which `DatetimeIndex` can parse, such as a quarterly string.

Previous behavior:

```
In [1]: pd.Timestamp('2012Q2')
Traceback
...
ValueError: Unable to parse 2012Q2

# Results in today's date.
In [2]: pd.Timestamp('2014')
Out [2]: 2014-08-12 00:00:00
```

v0.17.0 can parse them as below. It works on `DatetimeIndex` also.

New behavior:

```
In [63]: pd.Timestamp('2012Q2')
Out [63]: Timestamp('2012-04-01 00:00:00')

In [64]: pd.Timestamp('2014')
Out [64]: Timestamp('2014-01-01 00:00:00')

In [65]: pd.DatetimeIndex(['2012Q2', '2014'])
Out [65]: DatetimeIndex(['2012-04-01', '2014-01-01'], dtype='datetime64[ns]',
↪ freq=None)
```

Note: If you want to perform calculations based on today's date, use `Timestamp.now()` and `pandas.tseries.offsets`.

```
In [66]: import pandas.tseries.offsets as offsets

In [67]: pd.Timestamp.now()
```

(continues on next page)

(continued from previous page)

```
Out [67]: Timestamp('2020-06-17 17:53:27.046584')

In [68]: pd.Timestamp.now() + offsets.DateOffset(years=1)
Out [68]: Timestamp('2021-06-17 17:53:27.048514')
```

Changes to Index comparisons

Operator equal on Index should behavior similarly to Series ([GH9947](#), [GH10637](#))

Starting in v0.17.0, comparing Index objects of different lengths will raise a `ValueError`. This is to be consistent with the behavior of `Series`.

Previous behavior:

```
In [2]: pd.Index([1, 2, 3]) == pd.Index([1, 4, 5])
Out [2]: array([ True, False, False], dtype=bool)

In [3]: pd.Index([1, 2, 3]) == pd.Index([2])
Out [3]: array([False,  True, False], dtype=bool)

In [4]: pd.Index([1, 2, 3]) == pd.Index([1, 2])
Out [4]: False
```

New behavior:

```
In [8]: pd.Index([1, 2, 3]) == pd.Index([1, 4, 5])
Out [8]: array([ True, False, False], dtype=bool)

In [9]: pd.Index([1, 2, 3]) == pd.Index([2])
ValueError: Lengths must match to compare

In [10]: pd.Index([1, 2, 3]) == pd.Index([1, 2])
ValueError: Lengths must match to compare
```

Note that this is different from the `numpy` behavior where a comparison can be broadcast:

```
In [69]: np.array([1, 2, 3]) == np.array([1])
Out [69]: array([ True, False, False])
```

or it can return `False` if broadcasting can not be done:

```
In [70]: np.array([1, 2, 3]) == np.array([1, 2])
Out [70]: False
```

Changes to boolean comparisons vs. None

Boolean comparisons of a Series vs None will now be equivalent to comparing with `np.nan`, rather than raise `TypeError`. (GH1079).

```
In [71]: s = pd.Series(range(3))

In [72]: s.iloc[1] = None

In [73]: s
Out[73]:
0    0.0
1    NaN
2    2.0
Length: 3, dtype: float64
```

Previous behavior:

```
In [5]: s == None
TypeError: Could not compare <type 'NoneType'> type with Series
```

New behavior:

```
In [74]: s == None
Out[74]:
0    False
1    False
2    False
Length: 3, dtype: bool
```

Usually you simply want to know which values are null.

```
In [75]: s.isnull()
Out[75]:
0    False
1     True
2    False
Length: 3, dtype: bool
```

Warning: You generally will want to use `isnull/notnull` for these types of comparisons, as `isnull/notnull` tells you which elements are null. One has to be mindful that `nan`'s don't compare equal, but `None`'s do. Note that Pandas/numpy uses the fact that `np.nan != np.nan`, and treats `None` like `np.nan`.

```
In [76]: None == None
Out[76]: True

In [77]: np.nan == np.nan
Out[77]: False
```

HDFStore dropna behavior

The default behavior for HDFStore write functions with `format='table'` is now to keep rows that are all missing. Previously, the behavior was to drop rows that were all missing save the index. The previous behavior can be replicated using the `dropna=True` option. ([GH9382](#))

Previous behavior:

```
In [78]: df_with_missing = pd.DataFrame({'col1': [0, np.nan, 2],
....:                                   'col2': [1, np.nan, np.nan]})
....:

In [79]: df_with_missing
Out[79]:
   col1  col2
0    0.0    1.0
1    NaN    NaN
2    2.0    NaN

[3 rows x 2 columns]
```

```
In [27]:
df_with_missing.to_hdf('file.h5',
                       'df_with_missing',
                       format='table',
                       mode='w')

In [28]: pd.read_hdf('file.h5', 'df_with_missing')

Out [28]:
   col1  col2
0     0     1
2     2    NaN
```

New behavior:

```
In [80]: df_with_missing.to_hdf('file.h5',
....:                             'df_with_missing',
....:                             format='table',
....:                             mode='w')
....:

In [81]: pd.read_hdf('file.h5', 'df_with_missing')
Out[81]:
   col1  col2
0    0.0    1.0
1    NaN    NaN
2    2.0    NaN

[3 rows x 2 columns]
```

See the [docs](#) for more details.

Changes to `display.precision` option

The `display.precision` option has been clarified to refer to decimal places ([GH10451](#)).

Earlier versions of pandas would format floating point numbers to have one less decimal place than the value in `display.precision`.

```
In [1]: pd.set_option('display.precision', 2)

In [2]: pd.DataFrame({'x': [123.456789]})
Out[2]:
      x
0  123.5
```

If interpreting precision as “significant figures” this did work for scientific notation but that same interpretation did not work for values with standard formatting. It was also out of step with how numpy handles formatting.

Going forward the value of `display.precision` will directly control the number of places after the decimal, for regular formatting as well as scientific notation, similar to how numpy’s `precision` print option works.

```
In [82]: pd.set_option('display.precision', 2)

In [83]: pd.DataFrame({'x': [123.456789]})
Out[83]:
      x
0  123.46

[1 rows x 1 columns]
```

To preserve output behavior with prior versions the default value of `display.precision` has been reduced to 6 from 7.

Changes to `Categorical.unique`

`Categorical.unique` now returns new `Categoricals` with categories and codes that are unique, rather than returning `np.array` ([GH10508](#))

- unordered category: values and categories are sorted by appearance order.
- ordered category: values are sorted by appearance order, categories keep existing order.

```
In [84]: cat = pd.Categorical(['C', 'A', 'B', 'C'],
.....:                      categories=['A', 'B', 'C'],
.....:                      ordered=True)
.....:

In [85]: cat
Out[85]:
[C, A, B, C]
Categories (3, object): [A < B < C]

In [86]: cat.unique()
Out[86]:
[C, A, B]
Categories (3, object): [A < B < C]

In [87]: cat = pd.Categorical(['C', 'A', 'B', 'C'],
```

(continues on next page)

(continued from previous page)

```
.....:                                     categories=['A', 'B', 'C'])
.....:

In [88]: cat
Out[88]:
[C, A, B, C]
Categories (3, object): [A, B, C]

In [89]: cat.unique()
Out[89]:
[C, A, B]
Categories (3, object): [C, A, B]
```

Changes to `bool` passed as header in parsers

In earlier versions of pandas, if a `bool` was passed the `header` argument of `read_csv`, `read_excel`, or `read_html` it was implicitly converted to an integer, resulting in `header=0` for `False` and `header=1` for `True` ([GH6113](#))

A `bool` input to `header` will now raise a `TypeError`

```
In [29]: df = pd.read_csv('data.csv', header=False)
TypeError: Passing a bool to header is invalid. Use header=None for no header or
header=int or list-like of ints to specify the row(s) making up the column names
```

Other API changes

- Line and kde plot with `subplots=True` now uses default colors, not all black. Specify `color='k'` to draw all lines in black ([GH9894](#))
- Calling the `.value_counts()` method on a `Series` with a categorical dtype now returns a `Series` with a `CategoricalIndex` ([GH10704](#))
- The metadata properties of subclasses of pandas objects will now be serialized ([GH10553](#)).
- `groupby` using `Categorical` follows the same rule as `Categorical.unique` described above ([GH10508](#))
- When constructing `DataFrame` with an array of `complex64` dtype previously meant the corresponding column was automatically promoted to the `complex128` dtype. Pandas will now preserve the `itemsize` of the input for complex data ([GH10952](#))
- some numeric reduction operators would return `ValueError`, rather than `TypeError` on object types that includes strings and numbers ([GH11131](#))
- Passing currently unsupported `chunksize` argument to `read_excel` or `ExcelFile.parse` will now raise `NotImplementedError` ([GH8011](#))
- Allow an `ExcelFile` object to be passed into `read_excel` ([GH11198](#))
- `DatetimeIndex.union` does not infer `freq` if self and the input have `None` as `freq` ([GH11086](#))
- `NaT`'s methods now either raise `ValueError`, or return `np.nan` or `NaT` ([GH9513](#))

Behavior	Methods
<code>return np.nan</code>	<code>weekday, isoweekday</code>
<code>return NaT</code>	<code>date, now, replace, to_datetime, today</code>
<code>return np.datetime64('NaT')</code>	<code>to_datetime64</code> (unchanged)
<code>raise ValueError</code>	All other public methods (names not beginning with underscores)

Deprecations

- For `Series` the following indexing functions are deprecated ([GH10177](#)).

Deprecated Function	Replacement
<code>.irow(i)</code>	<code>.iloc[i]</code> or <code>.iat[i]</code>
<code>.iget(i)</code>	<code>.iloc[i]</code> or <code>.iat[i]</code>
<code>.iget_value(i)</code>	<code>.iloc[i]</code> or <code>.iat[i]</code>

- For `DataFrame` the following indexing functions are deprecated ([GH10177](#)).

Deprecated Function	Replacement
<code>.irow(i)</code>	<code>.iloc[i]</code>
<code>.iget_value(i, j)</code>	<code>.iloc[i, j]</code> or <code>.iat[i, j]</code>
<code>.icol(j)</code>	<code>.iloc[:, j]</code>

Note: These indexing function have been deprecated in the documentation since 0.11.0.

- `Categorical.name` was deprecated to make `Categorical` more `numpy.ndarray` like. Use `Series(cat, name="whatever")` instead ([GH10482](#)).
- Setting missing values (NaN) in a `Categorical`'s categories will issue a warning ([GH10748](#)). You can still have missing values in the values.
- `drop_duplicates` and `drop_duplicates`'s `take_last` keyword was deprecated in favor of `keep`. ([GH6511](#), [GH8505](#))
- `Series.nsmallest` and `nlargest`'s `take_last` keyword was deprecated in favor of `keep`. ([GH10792](#))
- `DataFrame.combineAdd` and `DataFrame.combineMult` are deprecated. They can easily be replaced by using the `add` and `mul` methods: `DataFrame.add(other, fill_value=0)` and `DataFrame.mul(other, fill_value=1.)` ([GH10735](#)).
- `TimeSeries` deprecated in favor of `Series` (note that this has been an alias since 0.13.0), ([GH10890](#))
- `SparsePanel` deprecated and will be removed in a future version ([GH11157](#)).
- `Series.is_time_series` deprecated in favor of `Series.index.is_all_dates` ([GH11135](#))
- Legacy offsets (like `'A@JAN'`) are deprecated (note that this has been alias since 0.8.0) ([GH10878](#))
- `WidePanel` deprecated in favor of `Panel`, `LongPanel` in favor of `DataFrame` (note these have been aliases since < 0.11.0), ([GH10892](#))
- `DataFrame.convert_objects` has been deprecated in favor of type-specific functions `pd.to_datetime`, `pd.to_timestamp` and `pd.to_numeric` (new in 0.17.0) ([GH11133](#)).

Removal of prior version deprecations/changes

- Removal of `na_last` parameters from `Series.order()` and `Series.sort()`, in favor of `na_position`. (GH5231)
- Remove of `percentile_width` from `.describe()`, in favor of `percentiles`. (GH7088)
- Removal of `colSpace` parameter from `DataFrame.to_string()`, in favor of `col_space`, circa 0.8.0 version.
- Removal of automatic time-series broadcasting (GH2304)

```
In [90]: np.random.seed(1234)

In [91]: df = pd.DataFrame(np.random.randn(5, 2),
.....:                    columns=list('AB'),
.....:                    index=pd.date_range('2013-01-01', periods=5))
.....:

In [92]: df
Out[92]:
```

	A	B
2013-01-01	0.471435	-1.190976
2013-01-02	1.432707	-0.312652
2013-01-03	-0.720589	0.887163
2013-01-04	0.859588	-0.636524
2013-01-05	0.015696	-2.242685

```
[5 rows x 2 columns]
```

Previously

```
In [3]: df + df.A
FutureWarning: TimeSeries broadcasting along DataFrame index by default is
↳ deprecated.
Please use DataFrame.<op> to explicitly broadcast arithmetic operations along the
↳ index

Out[3]:
```

	A	B
2013-01-01	0.942870	-0.719541
2013-01-02	2.865414	1.120055
2013-01-03	-1.441177	0.166574
2013-01-04	1.719177	0.223065
2013-01-05	0.031393	-2.226989

Current

```
In [93]: df.add(df.A, axis='index')
Out[93]:
```

	A	B
2013-01-01	0.942870	-0.719541
2013-01-02	2.865414	1.120055
2013-01-03	-1.441177	0.166574
2013-01-04	1.719177	0.223065
2013-01-05	0.031393	-2.226989

```
[5 rows x 2 columns]
```

- Remove `table` keyword in `HDFStore.put/append`, in favor of using `format=` (GH4645)
- Remove `kind` in `read_excel/ExcelFile` as its unused (GH4712)
- Remove `infer_type` keyword from `pd.read_html` as its unused (GH4770, GH7032)
- Remove `offset` and `timeRule` keywords from `Series.tshift/shift`, in favor of `freq` (GH4853, GH4864)
- Remove `pd.load/pd.save` aliases in favor of `pd.to_pickle/pd.read_pickle` (GH3787)

Performance improvements

- Development support for benchmarking with the `Air Speed Velocity` library (GH8361)
- Added `vbench` benchmarks for alternative `ExcelWriter` engines and reading Excel files (GH7171)
- Performance improvements in `Categorical.value_counts` (GH10804)
- Performance improvements in `SeriesGroupBy.nunique` and `SeriesGroupBy.value_counts` and `SeriesGroupby.transform` (GH10820, GH11077)
- Performance improvements in `DataFrame.drop_duplicates` with integer dtypes (GH10917)
- Performance improvements in `DataFrame.duplicated` with wide frames. (GH10161, GH11180)
- 4x improvement in `timedelta` string parsing (GH6755, GH10426)
- 8x improvement in `timedelta64` and `datetime64` ops (GH6755)
- Significantly improved performance of indexing `MultiIndex` with slicers (GH10287)
- 8x improvement in `iloc` using list-like input (GH10791)
- Improved performance of `Series.isin` for `datetimelike`/integer Series (GH10287)
- 20x improvement in `concat` of `Categoricals` when categories are identical (GH10587)
- Improved performance of `to_datetime` when specified format string is `ISO8601` (GH10178)
- 2x improvement of `Series.value_counts` for float dtype (GH10821)
- Enable `infer_datetime_format` in `to_datetime` when date components do not have 0 padding (GH11142)
- Regression from 0.16.1 in constructing `DataFrame` from nested dictionary (GH11084)
- Performance improvements in addition/subtraction operations for `DateOffset` with `Series` or `DatetimeIndex` (GH10744, GH11205)

Bug fixes

- Bug in incorrect computation of `.mean()` on `timedelta64[ns]` because of overflow (GH9442)
- Bug in `.isin` on older numpies (GH11232)
- Bug in `DataFrame.to_html(index=False)` renders unnecessary name row (GH10344)
- Bug in `DataFrame.to_latex()` the `column_format` argument could not be passed (GH9402)
- Bug in `DatetimeIndex` when localizing with `NaT` (GH10477)
- Bug in `Series.dt` ops in preserving meta-data (GH10477)
- Bug in preserving `NaT` when passed in an otherwise invalid `to_datetime` construction (GH10477)

- Bug in `DataFrame.apply` when function returns categorical series. (GH9573)
- Bug in `to_datetime` with invalid dates and formats supplied (GH10154)
- Bug in `Index.drop_duplicates` dropping name(s) (GH10115)
- Bug in `Series.quantile` dropping name (GH10881)
- Bug in `pd.Series` when setting a value on an empty `Series` whose index has a frequency. (GH10193)
- Bug in `pd.Series.interpolate` with invalid order keyword values. (GH10633)
- Bug in `DataFrame.plot` raises `ValueError` when color name is specified by multiple characters (GH10387)
- Bug in `Index` construction with a mixed list of tuples (GH10697)
- Bug in `DataFrame.reset_index` when index contains `NaT`. (GH10388)
- Bug in `ExcelReader` when worksheet is empty (GH6403)
- Bug in `BinGrouper.group_info` where returned values are not compatible with base class (GH10914)
- Bug in clearing the cache on `DataFrame.pop` and a subsequent inplace op (GH10912)
- Bug in indexing with a mixed-integer `Index` causing an `ImportError` (GH10610)
- Bug in `Series.count` when index has nulls (GH10946)
- Bug in pickling of a non-regular freq `DatetimeIndex` (GH11002)
- Bug causing `DataFrame.where` to not respect the `axis` parameter when the frame has a symmetric shape. (GH9736)
- Bug in `Table.select_column` where name is not preserved (GH10392)
- Bug in `offsets.generate_range` where start and end have finer precision than offset (GH9907)
- Bug in `pd.rolling_*` where `Series.name` would be lost in the output (GH10565)
- Bug in `stack` when index or columns are not unique. (GH10417)
- Bug in setting a `Panel` when an axis has a `MultiIndex` (GH10360)
- Bug in `USFederalHolidayCalendar` where `USMemorialDay` and `USMartinLutherKingJr` were incorrect (GH10278 and GH9760)
- Bug in `.sample()` where returned object, if set, gives unnecessary `SettingWithCopyWarning` (GH10738)
- Bug in `.sample()` where weights passed as `Series` were not aligned along axis before being treated positionally, potentially causing problems if weight indices were not aligned with sampled object. (GH10738)
- Regression fixed in (GH9311, GH6620, GH9345), where `groupby` with a datetime-like converting to float with certain aggregators (GH10979)
- Bug in `DataFrame.interpolate` with `axis=1` and `inplace=True` (GH10395)
- Bug in `io.sql.get_schema` when specifying multiple columns as primary key (GH10385).
- Bug in `groupby(sort=False)` with datetime-like `Categorical` raises `ValueError` (GH10505)
- Bug in `groupby(axis=1)` with `filter()` throws `IndexError` (GH11041)
- Bug in `test_categorical` on big-endian builds (GH10425)
- Bug in `Series.shift` and `DataFrame.shift` not supporting categorical data (GH9416)
- Bug in `Series.map` using categorical `Series` raises `AttributeError` (GH10324)