




**All checks have passed**[Hide all checks](#)

2 successful checks

 **continuous-integration/travis-ci/pr** — The Travis CI build passed[Details](#)

 **pandas-dev.pandas** Successful in 36m — Build #20190109.23 succeeded[Details](#)

**This branch has no conflicts with the base branch**

Merging can be performed automatically.

Squash and merge

 or view [command line instructions](#).

Note: Each time you push to *your* fork, a *new* run of the tests will be triggered on the CI. You can enable the auto-cancel feature, which removes any non-currently-running tests for that same pull-request, for [Travis-CI here](#).

Test-driven development/code writing

pandas is serious about testing and strongly encourages contributors to embrace [test-driven development \(TDD\)](#). This development process “relies on the repetition of a very short development cycle: first the developer writes an (initially failing) automated test case that defines a desired improvement or new function, then produces the minimum amount of code to pass that test.” So, before actually writing any code, you should write your tests. Often the test can be taken from the original GitHub issue. However, it is always worth considering additional use cases and writing corresponding tests.

Adding tests is one of the most common requests after code is pushed to *pandas*. Therefore, it is worth getting in the habit of writing tests ahead of time so this is never an issue.

Like many packages, *pandas* uses [pytest](#) and the convenient extensions in [numpy.testing](#).

Note: The earliest supported pytest version is 5.0.1.

Writing tests

All tests should go into the `tests` subdirectory of the specific package. This folder contains many current examples of tests, and we suggest looking to these for inspiration. If your test requires working with files or network connectivity, there is more information on the [testing page](#) of the wiki.

The `pandas._testing` module has many special `assert` functions that make it easier to make statements about whether `Series` or `DataFrame` objects are equivalent. The easiest way to verify that your code is correct is to explicitly construct the result you expect, then compare the actual result to the expected correct result:

```
def test_pivot(self):
    data = {
        'index' : ['A', 'B', 'C', 'C', 'B', 'A'],
        'columns' : ['One', 'One', 'One', 'Two', 'Two', 'Two'],
        'values' : [1., 2., 3., 3., 2., 1.]
```

(continues on next page)

(continued from previous page)

```

}

frame = DataFrame(data)
pivoted = frame.pivot(index='index', columns='columns', values='values')

expected = DataFrame({
    'One' : {'A' : 1., 'B' : 2., 'C' : 3.},
    'Two' : {'A' : 1., 'B' : 2., 'C' : 3.}
})

assert_frame_equal(pivoted, expected)

```

Please remember to add the Github Issue Number as a comment to a new test. E.g. “# brief comment, see GH#28907”

Transitioning to pytest

pandas existing test structure is *mostly* class-based, meaning that you will typically find tests wrapped in a class.

```

class TestReallyCoolFeature:
    pass

```

Going forward, we are moving to a more *functional* style using the `pytest` framework, which offers a richer testing framework that will facilitate testing and developing. Thus, instead of writing test classes, we will write test functions like this:

```

def test_really_cool_feature():
    pass

```

Using pytest

Here is an example of a self-contained set of tests that illustrate multiple features that we like to use.

- functional style: tests are like `test_*` and *only* take arguments that are either fixtures or parameters
- `pytest.mark` can be used to set metadata on test functions, e.g. `skip` or `xfail`.
- using `parametrize`: allow testing of multiple cases
- to set a mark on a parameter, `pytest.param(..., marks=...)` syntax should be used
- `fixture`, code for object construction, on a per-test basis
- using bare `assert` for scalars and truth-testing
- `tm.assert_series_equal` (and its counter part `tm.assert_frame_equal`), for pandas object comparisons.
- the typical pattern of constructing an `expected` and comparing versus the `result`

We would name this file `test_cool_feature.py` and put in an appropriate place in the `pandas/tests/` structure.

```

import pytest
import numpy as np
import pandas as pd

```

(continues on next page)

(continued from previous page)

```

@pytest.mark.parametrize('dtype', ['int8', 'int16', 'int32', 'int64'])
def test_dtypes(dtype):
    assert str(np.dtype(dtype)) == dtype

@pytest.mark.parametrize(
    'dtype', ['float32', pytest.param('int16', marks=pytest.mark.skip),
        pytest.param('int32', marks=pytest.mark.xfail(
            reason='to show how it works'))])
def test_mark(dtype):
    assert str(np.dtype(dtype)) == 'float32'

@pytest.fixture
def series():
    return pd.Series([1, 2, 3])

@pytest.fixture(params=['int8', 'int16', 'int32', 'int64'])
def dtype(request):
    return request.param

def test_series(series, dtype):
    result = series.astype(dtype)
    assert result.dtype == dtype

    expected = pd.Series([1, 2, 3], dtype=dtype)
    tm.assert_series_equal(result, expected)

```

A test run of this yields

```

((pandas) bash-3.2$ pytest test_cool_feature.py -v
===== test session starts =====
platform darwin -- Python 3.6.2, pytest-3.6.0, py-1.4.31, pluggy-0.4.0
collected 11 items

tester.py::test_dtypes[int8] PASSED
tester.py::test_dtypes[int16] PASSED
tester.py::test_dtypes[int32] PASSED
tester.py::test_dtypes[int64] PASSED
tester.py::test_mark[float32] PASSED
tester.py::test_mark[int16] SKIPPED
tester.py::test_mark[int32] xfail
tester.py::test_series[int8] PASSED
tester.py::test_series[int16] PASSED
tester.py::test_series[int32] PASSED
tester.py::test_series[int64] PASSED

```

Tests that we have parametrized are now accessible via the test name, for example we could run these with `-k int8` to sub-select *only* those tests which match `int8`.

```

((pandas) bash-3.2$ pytest test_cool_feature.py -v -k int8
===== test session starts =====
platform darwin -- Python 3.6.2, pytest-3.6.0, py-1.4.31, pluggy-0.4.0

```

(continues on next page)

(continued from previous page)

```
collected 11 items

test_cool_feature.py::test_dtypes[int8] PASSED
test_cool_feature.py::test_series[int8] PASSED
```

Using hypothesis

Hypothesis is a library for property-based testing. Instead of explicitly parametrizing a test, you can describe *all* valid inputs and let Hypothesis try to find a failing input. Even better, no matter how many random examples it tries, Hypothesis always reports a single minimal counterexample to your assertions - often an example that you would never have thought to test.

See [Getting Started with Hypothesis](#) for more of an introduction, then [refer to the Hypothesis documentation](#) for details.

```
import json
from hypothesis import given, strategies as st

any_json_value = st.deferred(lambda: st.one_of(
    st.none(), st.booleans(), st.floats(allow_nan=False), st.text(),
    st.lists(any_json_value), st.dictionaries(st.text(), any_json_value)
))

@given(value=any_json_value)
def test_json_roundtrip(value):
    result = json.loads(json.dumps(value))
    assert value == result
```

This test shows off several useful features of Hypothesis, as well as demonstrating a good use-case: checking properties that should hold over a large or complicated domain of inputs.

To keep the Pandas test suite running quickly, parametrized tests are preferred if the inputs or logic are simple, with Hypothesis tests reserved for cases with complex logic or where there are too many combinations of options or subtle interactions to test (or think of!) all of them.

Testing warnings

By default, one of pandas CI workers will fail if any unhandled warnings are emitted.

If your change involves checking that a warning is actually emitted, use `tm.assert_produces_warning(ExpectedWarning)`.

```
import pandas._testing as tm

df = pd.DataFrame()
with tm.assert_produces_warning(FutureWarning):
    df.some_operation()
```

We prefer this to the `pytest.warns` context manager because ours checks that the warning's `stacklevel` is set correctly. The `stacklevel` is what ensure the *user's* file name and line number is printed in the warning, rather than something internal to pandas. It represents the number of function calls from user code (e.g. `df.some_operation()`) to the function that actually emits the warning. Our linter will fail the build if you use `pytest.warns` in a test.

If you have a test that would emit a warning, but you aren't actually testing the warning itself (say because it's going to be removed in the future, or because we're matching a 3rd-party library's behavior), then use `pytest.mark.filterwarnings` to ignore the error.

```
@pytest.mark.filterwarnings("ignore:msg:category")
def test_thing(self):
    ...
```

If the test generates a warning of class `category` whose message starts with `msg`, the warning will be ignored and the test will pass.

If you need finer-grained control, you can use Python's usual `warnings` module to control whether a warning is ignored / raised at different places within a single test.

```
with warnings.catch_warnings():
    warnings.simplefilter("ignore", FutureWarning)
    # Or use warnings.filterwarnings(...)
```

Alternatively, consider breaking up the unit test.

Running the test suite

The tests can then be run directly inside your Git clone (without having to install *pandas*) by typing:

```
pytest pandas
```

The tests suite is exhaustive and takes around 20 minutes to run. Often it is worth running only a subset of tests first around your changes before running the entire suite.

The easiest way to do this is with:

```
pytest pandas/path/to/test.py -k regex_matching_test_name
```

Or with one of the following constructs:

```
pytest pandas/tests/[test-module].py
pytest pandas/tests/[test-module].py::[TestClass]
pytest pandas/tests/[test-module].py::[TestClass]::[test_method]
```

Using `pytest-xdist`, one can speed up local testing on multicore machines. To use this feature, you will need to install *pytest-xdist* via:

```
pip install pytest-xdist
```

Two scripts are provided to assist with this. These scripts distribute testing across 4 threads.

On Unix variants, one can type:

```
test_fast.sh
```

On Windows, one can type:

```
test_fast.bat
```

This can significantly reduce the time it takes to locally run tests before submitting a pull request.

For more, see the [pytest](#) documentation.

Furthermore one can run

```
pd.test()
```

with an imported pandas to run tests similarly.

Running the performance test suite

Performance matters and it is worth considering whether your code has introduced performance regressions. *pandas* is in the process of migrating to [asv benchmarks](#) to enable easy monitoring of the performance of critical *pandas* operations. These benchmarks are all found in the `pandas/asv_bench` directory. *asv* supports both python2 and python3.

To use all features of *asv*, you will need either `conda` or `virtualenv`. For more details please check the [asv installation webpage](#).

To install *asv*:

```
pip install git+https://github.com/spacetelescope/asv
```

If you need to run a benchmark, change your directory to `asv_bench/` and run:

```
asv continuous -f 1.1 upstream/master HEAD
```

You can replace `HEAD` with the name of the branch you are working on, and report benchmarks that changed by more than 10%. The command uses `conda` by default for creating the benchmark environments. If you want to use `virtualenv` instead, write:

```
asv continuous -f 1.1 -E virtualenv upstream/master HEAD
```

The `-E virtualenv` option should be added to all *asv* commands that run benchmarks. The default value is defined in `asv.conf.json`.

Running the full test suite can take up to one hour and use up to 3GB of RAM. Usually it is sufficient to paste only a subset of the results into the pull request to show that the committed changes do not cause unexpected performance regressions. You can run specific benchmarks using the `-b` flag, which takes a regular expression. For example, this will only run tests from a `pandas/asv_bench/benchmarks/groupby.py` file:

```
asv continuous -f 1.1 upstream/master HEAD -b ^groupby
```

If you want to only run a specific group of tests from a file, you can do it using `.` as a separator. For example:

```
asv continuous -f 1.1 upstream/master HEAD -b groupby.GroupByMethods
```

will only run the `GroupByMethods` benchmark defined in `groupby.py`.

You can also run the benchmark suite using the version of *pandas* already installed in your current Python environment. This can be useful if you do not have `virtualenv` or `conda`, or are using the `setup.py develop` approach discussed above; for the in-place build you need to set `PYTHONPATH`, e.g. `PYTHONPATH="$PWD/.."` *asv* [remaining arguments]. You can run benchmarks using an existing Python environment by:

```
asv run -e -E existing
```

or, to use a specific Python interpreter,:

```
asv run -e -E existing:python3.6
```

This will display `stderr` from the benchmarks, and use your local `python` that comes from your `$PATH`.

Information on how to write a benchmark and how to use *asv* can be found in the [asv documentation](#).

Documenting your code

Changes should be reflected in the release notes located in `doc/source/whatsnew/vx.y.z.rst`. This file contains an ongoing change log for each release. Add an entry to this file to document your fix, enhancement or (unavoidable) breaking change. Make sure to include the GitHub issue number when adding your entry (using `:issue:`1234`` where 1234 is the issue/pull request number).

If your code is an enhancement, it is most likely necessary to add usage examples to the existing documentation. This can be done following the section regarding documentation [above](#). Further, to let users know when this feature was added, the `versionadded` directive is used. The sphinx syntax for that is:

```
.. versionadded:: 0.21.0
```

This will put the text *New in version 0.21.0* wherever you put the sphinx directive. This should also be put in the docstring when adding a new function or method ([example](#)) or a new keyword argument ([example](#)).

4.1.6 Contributing your changes to *pandas*

Committing your code

Keep style fixes to a separate commit to make your pull request more readable.

Once you've made changes, you can see them by typing:

```
git status
```

If you have created a new file, it is not being tracked by git. Add it by typing:

```
git add path/to/file-to-be-added.py
```

Doing 'git status' again should give something like:

```
# On branch shiny-new-feature
#
#       modified:   /relative/path/to/file-you-added.py
#
```

Finally, commit your changes to your local repository with an explanatory message. *Pandas* uses a convention for commit message prefixes and layout. Here are some common prefixes along with general guidelines for when to use them:

- ENH: Enhancement, new functionality
- BUG: Bug fix
- DOC: Additions/updates to documentation
- TST: Additions/updates to tests
- BLD: Updates to the build process/scripts
- PERF: Performance improvement
- TYP: Type annotations
- CLN: Code cleanup

The following defines how a commit message should be structured. Please reference the relevant GitHub issues in your commit message using GH1234 or #1234. Either style is fine, but the former is generally preferred:

- a subject line with < 80 chars.
- One blank line.
- Optionally, a commit message body.

Now you can commit your changes in your local repository:

```
git commit -m
```

Pushing your changes

When you want your changes to appear publicly on your GitHub page, push your forked feature branch's commits:

```
git push origin shiny-new-feature
```

Here `origin` is the default name given to your remote repository on GitHub. You can see the remote repositories:

```
git remote -v
```

If you added the upstream repository as described above you will see something like:

```
origin  git@github.com:yourname/pandas.git (fetch)
origin  git@github.com:yourname/pandas.git (push)
upstream      git://github.com/pandas-dev/pandas.git (fetch)
upstream      git://github.com/pandas-dev/pandas.git (push)
```

Now your code is on GitHub, but it is not yet a part of the *pandas* project. For that to happen, a pull request needs to be submitted on GitHub.

Review your code

When you're ready to ask for a code review, file a pull request. Before you do, once again make sure that you have followed all the guidelines outlined in this document regarding code style, tests, performance tests, and documentation. You should also double check your branch changes against the branch it was based on:

1. Navigate to your repository on GitHub – <https://github.com/your-user-name/pandas>
2. Click on `Branches`
3. Click on the `Compare` button for your feature branch
4. Select the `base` and `compare` branches, if necessary. This will be `master` and `shiny-new-feature`, respectively.

Finally, make the pull request

If everything looks good, you are ready to make a pull request. A pull request is how code from a local repository becomes available to the GitHub community and can be looked at and eventually merged into the master version. This pull request and its associated changes will eventually be committed to the master branch and available in the next release. To submit a pull request:

1. Navigate to your repository on GitHub
2. Click on the `Pull Request` button
3. You can then click on `Commits` and `Files Changed` to make sure everything looks okay one last time

4. Write a description of your changes in the `Preview Discussion` tab
5. Click `Send Pull Request`.

This request then goes to the repository maintainers, and they will review the code.

Updating your pull request

Based on the review you get on your pull request, you will probably need to make some changes to the code. In that case, you can make them in your branch, add a new commit to that branch, push it to GitHub, and the pull request will be automatically updated. Pushing them to GitHub again is done by:

```
git push origin shiny-new-feature
```

This will automatically update your pull request with the latest code and restart the *Continuous Integration* tests.

Another reason you might need to update your pull request is to solve conflicts with changes that have been merged into the master branch since you opened your pull request.

To do this, you need to “merge upstream master” in your branch:

```
git checkout shiny-new-feature
git fetch upstream
git merge upstream/master
```

If there are no conflicts (or they could be fixed automatically), a file with a default commit message will open, and you can simply save and quit this file.

If there are merge conflicts, you need to solve those conflicts. See for example at <https://help.github.com/articles/resolving-a-merge-conflict-using-the-command-line/> for an explanation on how to do this. Once the conflicts are merged and the files where the conflicts were solved are added, you can run `git commit` to save those fixes.

If you have uncommitted changes at the moment you want to update the branch with master, you will need to `stash` them prior to updating (see the [stash docs](#)). This will effectively store your changes and they can be reapplied after updating.

After the feature branch has been update locally, you can now update your pull request by pushing to the branch on GitHub:

```
git push origin shiny-new-feature
```

Delete your merged branch (optional)

Once your feature branch is accepted into upstream, you’ll probably want to get rid of the branch. First, merge upstream master into your branch so git knows it is safe to delete your branch:

```
git fetch upstream
git checkout master
git merge upstream/master
```

Then you can do:

```
git branch -d shiny-new-feature
```

Make sure you use a lower-case `-d`, or else git won’t warn you if your feature branch has not actually been merged.

The branch will still exist on GitHub, so to delete it there do:

```
git push origin --delete shiny-new-feature
```

4.2 pandas code style guide

Table of contents:

- *Patterns*
 - *foo.__class__*
- *String formatting*
 - *Concatenated strings*
 - * *f-strings*
 - * *White spaces*
 - *Representation function (aka 'repr()')*

4.2.1 Patterns

foo.__class__

pandas uses 'type(foo)' instead 'foo.__class__' as it is making the code more readable.

For example:

Good:

```
foo = "bar"
type(foo)
```

Bad:

```
foo = "bar"
foo.__class__
```

4.2.2 String formatting

Concatenated strings

f-strings

pandas uses f-strings formatting instead of '%' and '.format()' string formatters.

The convention of using f-strings on a string that is concatenated over several lines, is to prefix only the lines containing the value needs to be interpreted.

For example:

Good:

```
foo = "old_function"
bar = "new_function"

my_warning_message = (
    f"Warning, {foo} is deprecated, "
    "please use the new and way better "
    f"{bar}"
)
```

Bad:

```
foo = "old_function"
bar = "new_function"

my_warning_message = (
    f"Warning, {foo} is deprecated, "
    f"please use the new and way better "
    f"{bar}"
)
```

White spaces

Putting the white space only at the end of the previous line, so there is no whitespace at the beginning of the concatenated string.

For example:

Good:

```
example_string = (
    "Some long concatenated string, "
    "with good placement of the "
    "whitespaces"
)
```

Bad:

```
example_string = (
    "Some long concatenated string,"
    " with bad placement of the"
    " whitespaces"
)
```

Representation function (aka 'repr()')

pandas uses 'repr()' instead of '%r' and '!r'.

The use of 'repr()' will only happen when the value is not an obvious string.

For example:

Good:

```
value = str
f"Unknown recived value, got: {repr(value)}"
```

Good:

```
value = str
f"Unknown received type, got: '{type(value).__name__}'"
```

4.3 Pandas Maintenance

This guide is for pandas' maintainers. It may also be interesting to contributors looking to understand the pandas development process and what steps are necessary to become a maintainer.

The main contributing guide is available at [Contributing to pandas](#).

4.3.1 Roles

Pandas uses two levels of permissions: **triage** and **core** team members.

Triage members can label and close issues and pull requests.

Core team members can label and close issues and pull request, and can merge pull requests.

GitHub publishes the full [list of permissions](#).

4.3.2 Tasks

Pandas is largely a volunteer project, so these tasks shouldn't be read as "expectations" of triage and maintainers. Rather, they're general descriptions of what it means to be a maintainer.

- Triage newly filed issues (see [Issue Triage](#))
- Review newly opened pull requests
- Respond to updates on existing issues and pull requests
- Drive discussion and decisions on stalled issues and pull requests
- Provide experience / wisdom on API design questions to ensure consistency and maintainability
- Project organization (run / attend developer meetings, represent pandas)

<http://matthewrocklin.com/blog/2019/05/18/maintainer> may be interesting background reading.

4.3.3 Issue Triage

Here's a typical workflow for triaging a newly opened issue.

1. Thank the reporter for opening an issue

The issue tracker is many people's first interaction with the pandas project itself, beyond just using the library. As such, we want it to be a welcoming, pleasant experience.

2. Is the necessary information provided?

Ideally reporters would fill out the issue template, but many don't. If crucial information (like the version of pandas they used), is missing feel free to ask for that and label the issue with "Needs info". The report should follow the guidelines in [Bug reports and enhancement requests](#). You may want to link to that if they didn't follow the template.

Make sure that the title accurately reflects the issue. Edit it yourself if it's not clear.

3. Is this a duplicate issue?

We have many open issues. If a new issue is clearly a duplicate, label the new issue as “Duplicate” assign the milestone “No Action”, and close the issue with a link to the original issue. Make sure to still thank the reporter, and encourage them to chime in on the original issue, and perhaps try to fix it.

If the new issue provides relevant information, such as a better or slightly different example, add it to the original issue as a comment or an edit to the original post.

4. Is the issue minimal and reproducible?

For bug reports, we ask that the reporter provide a minimal reproducible example. See <http://matthewrocklin.com/blog/work/2018/02/28/minimal-bug-reports> for a good explanation. If the example is not reproducible, or if it's *clearly* not minimal, feel free to ask the reporter if they can provide an example or simplify the provided one. Do acknowledge that writing minimal reproducible examples is hard work. If the reporter is struggling, you can try to write one yourself and we'll edit the original post to include it.

If a reproducible example can't be provided, add the “Needs info” label.

If a reproducible example is provided, but you see a simplification, edit the original post with your simpler reproducible example.

5. Is this a clearly defined feature request?

Generally, pandas prefers to discuss and design new features in issues, before a pull request is made. Encourage the submitter to include a proposed API for the new feature. Having them write a full docstring is a good way to pin down specifics.

We'll need a discussion from several pandas maintainers before deciding whether the proposal is in scope for pandas.

6. Is this a usage question?

We prefer that usage questions are asked on StackOverflow with the pandas tag. <https://stackoverflow.com/questions/tagged/pandas>

If it's easy to answer, feel free to link to the relevant documentation section, let them know that in the future this kind of question should be on StackOverflow, and close the issue.

7. What labels and milestones should I add?

Apply the relevant labels. This is a bit of an art, and comes with experience. Look at similar issues to get a feel for how things are labeled.

If the issue is clearly defined and the fix seems relatively straightforward, label the issue as “Good first issue”.

Typically, new issues will be assigned the “Contributions welcome” milestone, unless it's known that this issue should be addressed in a specific release (say because it's a large regression).

4.3.4 Closing Issues

Be delicate here: many people interpret closing an issue as us saying that the conversation is over. It's typically best to give the reporter some time to respond or self-close their issue if it's determined that the behavior is not a bug, or the feature is out of scope. Sometimes reporters just go away though, and we'll close the issue after the conversation has died.

4.3.5 Reviewing Pull Requests

Anybody can review a pull request: regular contributors, triagers, or core-team members. Here are some guidelines to check.

- Tests should be in a sensible location.
- New public APIs should be included somewhere in `doc/source/reference/`.
- New / changed API should use the `versionadded` or `versionchanged` directives in the docstring.
- User-facing changes should have a whatsnew in the appropriate file.
- Regression tests should reference the original GitHub issue number like `# GH-1234`.

4.3.6 Cleaning up old Issues

Every open issue in pandas has a cost. Open issues make finding duplicates harder, and can make it harder to know what needs to be done in pandas. That said, closing issues isn't a goal on its own. Our goal is to make pandas the best it can be, and that's best done by ensuring that the quality of our open issues is high.

Occasionally, bugs are fixed but the issue isn't linked to in the Pull Request. In these cases, comment that "This has been fixed, but could use a test." and label the issue as "Good First Issue" and "Needs Test".

If an older issue doesn't follow our issue template, edit the original post to include a minimal example, the actual output, and the expected output. Uniformity in issue reports is valuable.

If an older issue lacks a reproducible example, label it as "Needs Info" and ask them to provide one (or write one yourself if possible). If one isn't provide reasonably soon, close it according to the policies in *Closing Issues*.

4.3.7 Cleaning up old Pull Requests

Occasionally, contributors are unable to finish off a pull request. If some time has passed (two weeks, say) since the last review requesting changes, gently ask if they're still interested in working on this. If another two weeks or so passes with no response, thank them for their work and close the pull request. Comment on the original issue that "There's a stalled PR at #1234 that may be helpful.", and perhaps label the issue as "Good first issue" if the PR was relatively close to being accepted.

Additionally, core-team members can push to contributors branches. This can be helpful for pushing an important PR across the line, or for fixing a small merge conflict.

4.3.8 Becoming a pandas maintainer

The full process is outlined in our [governance documents](#). In summary, we're happy to give triage permissions to anyone who shows interest by being helpful on the issue tracker.

The current list of core-team members is at <https://github.com/pandas-dev/pandas-governance/blob/master/people.md>

4.4 Internals

This section will provide a look into some of pandas internals. It's primarily intended for developers of pandas itself.

4.4.1 Indexing

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
- `Float64Index`: a version of `Index` highly optimized for 64-bit float data
- `MultiIndex`: the standard hierarchical index object
- `DatetimeIndex`: An `Index` object with `Timestamp` boxed elements (impl are the int64 values)
- `TimedeltaIndex`: An `Index` object with `Timedelta` boxed elements (impl are the in64 values)
- `PeriodIndex`: An `Index` object with `Period` elements

There are functions that make the creation of a regular index easy:

- `date_range`: fixed frequency date range generated from a time rule or `DateOffset`. An ndarray of Python datetime objects
- `period_range`: fixed frequency date range generated from a time rule or `DateOffset`. An ndarray of `Period` objects, representing timespans

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.

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
- `get_indexer_non_unique`: Computes the indexing vector for reindexing / data alignment purposes when the index is non-unique. 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`

MultIndex

Internally, the `MultiIndex` consists of a few things: the **levels**, the integer **codes** (until version 0.24 named *labels*), and the level **names**:

```
In [1]: index = pd.MultiIndex.from_product([range(3), ['one', 'two']],
...:                                     names=['first', 'second'])
...:
...:

In [2]: index
Out[2]:
MultiIndex([(0, 'one'),
            (0, 'two'),
            (1, 'one'),
            (1, 'two'),
            (2, 'one'),
            (2, 'two')],
           names=['first', 'second'])

In [3]: index.levels
Out[3]: FrozenList([[0, 1, 2], ['one', 'two']])

In [4]: index.codes
Out[4]: FrozenList([[0, 0, 1, 1, 2, 2], [0, 1, 0, 1, 0, 1]])

In [5]: index.names
Out[5]: FrozenList(['first', 'second'])
```

You can probably guess that the codes 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 codes 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 codes yourself, please be careful.

Values

Pandas extends NumPy's type system with custom types, like `Categorical` or datetimes with a timezone, so we have multiple notions of "values". For 1-D containers (Index classes and `Series`) we have the following convention:

- `cls._ndarray_values` is *always* a NumPy `ndarray`. Ideally, `_ndarray_values` is cheap to compute. For example, for a `Categorical`, this returns the codes, not the array of objects.
- `cls._values` refers to the "best possible" array. This could be an `ndarray`, `ExtensionArray`, or in Index subclass (note: we're in the process of removing the index subclasses here so that it's always an `ndarray` or `ExtensionArray`).

So, for example, `Series[category]._values` is a `Categorical`, while `Series[category]._ndarray_values` is the underlying codes.

4.4.2 Subclassing pandas data structures

This section has been moved to *Subclassing pandas data structures*.

4.5 Extending pandas

While pandas provides a rich set of methods, containers, and data types, your needs may not be fully satisfied. Pandas offers a few options for extending pandas.

4.5.1 Registering custom accessors

Libraries can use the decorators `pandas.api.extensions.register_dataframe_accessor()`, `pandas.api.extensions.register_series_accessor()`, and `pandas.api.extensions.register_index_accessor()`, to add additional “namespaces” to pandas objects. All of these follow a similar convention: you decorate a class, providing the name of attribute to add. The class’s `__init__` method gets the object being decorated. For example:

```
@pd.api.extensions.register_dataframe_accessor("geo")
class GeoAccessor:
    def __init__(self, pandas_obj):
        self._validate(pandas_obj)
        self._obj = pandas_obj

    @staticmethod
    def _validate(obj):
        # verify there is a column latitude and a column longitude
        if 'latitude' not in obj.columns or 'longitude' not in obj.columns:
            raise AttributeError("Must have 'latitude' and 'longitude'.")

    @property
    def center(self):
        # return the geographic center point of this DataFrame
        lat = self._obj.latitude
        lon = self._obj.longitude
        return (float(lon.mean()), float(lat.mean()))

    def plot(self):
        # plot this array's data on a map, e.g., using Cartopy
        pass
```

Now users can access your methods using the `geo` namespace:

```
>>> ds = pd.DataFrame({'longitude': np.linspace(0, 10),
...                    'latitude': np.linspace(0, 20)})
>>> ds.geo.center
(5.0, 10.0)
>>> ds.geo.plot()
# plots data on a map
```

This can be a convenient way to extend pandas objects without subclassing them. If you write a custom accessor, make a pull request adding it to our ecosystem page.

We highly recommend validating the data in your accessor’s `__init__`. In our `GeoAccessor`, we validate that the data contains the expected columns, raising an `AttributeError` when the validation fails. For a `Series` accessor, you should validate the `dtype` if the accessor applies only to certain `dtypes`.

4.5.2 Extension types

New in version 0.23.0.

Warning: The `pandas.api.extensions.ExtensionDtype` and `pandas.api.extensions.ExtensionArray` APIs are new and experimental. They may change between versions without warning.

Pandas defines an interface for implementing data types and arrays that *extend* NumPy's type system. Pandas itself uses the extension system for some types that aren't built into NumPy (categorical, period, interval, datetime with timezone).

Libraries can define a custom array and data type. When pandas encounters these objects, they will be handled properly (i.e. not converted to an ndarray of objects). Many methods like `pandas.isna()` will dispatch to the extension type's implementation.

If you're building a library that implements the interface, please publicize it on ecosystem.extensions.

The interface consists of two classes.

ExtensionDtype

A `pandas.api.extensions.ExtensionDtype` is similar to a `numpy.dtype` object. It describes the data type. Implementors are responsible for a few unique items like the name.

One particularly important item is the `type` property. This should be the class that is the scalar type for your data. For example, if you were writing an extension array for IP Address data, this might be `ipaddress.IPv4Address`.

See the [extension dtype source](#) for interface definition.

New in version 0.24.0.

`pandas.api.extension.ExtensionDtype` can be registered to pandas to allow creation via a string dtype name. This allows one to instantiate `Series` and `.astype()` with a registered string name, for example `'category'` is a registered string accessor for the `CategoricalDtype`.

See the [extension dtype dtypes](#) for more on how to register dtypes.

ExtensionArray

This class provides all the array-like functionality. ExtensionArrays are limited to 1 dimension. An ExtensionArray is linked to an ExtensionDtype via the `dtype` attribute.

Pandas makes no restrictions on how an extension array is created via its `__new__` or `__init__`, and puts no restrictions on how you store your data. We do require that your array be convertible to a NumPy array, even if this is relatively expensive (as it is for `Categorical`).

They may be backed by none, one, or many NumPy arrays. For example, `pandas.Categorical` is an extension array backed by two arrays, one for codes and one for categories. An array of IPv6 addresses may be backed by a NumPy structured array with two fields, one for the lower 64 bits and one for the upper 64 bits. Or they may be backed by some other storage type, like Python lists.

See the [extension array source](#) for the interface definition. The docstrings and comments contain guidance for properly implementing the interface.

ExtensionArray Operator Support

New in version 0.24.0.

By default, there are no operators defined for the class `ExtensionArray`. There are two approaches for providing operator support for your `ExtensionArray`:

1. Define each of the operators on your `ExtensionArray` subclass.
2. Use an operator implementation from pandas that depends on operators that are already defined on the underlying elements (scalars) of the `ExtensionArray`.

Note: Regardless of the approach, you may want to set `__array_priority__` if you want your implementation to be called when involved in binary operations with NumPy arrays.

For the first approach, you define selected operators, e.g., `__add__`, `__le__`, etc. that you want your `ExtensionArray` subclass to support.

The second approach assumes that the underlying elements (i.e., scalar type) of the `ExtensionArray` have the individual operators already defined. In other words, if your `ExtensionArray` named `MyExtensionArray` is implemented so that each element is an instance of the class `MyExtensionElement`, then if the operators are defined for `MyExtensionElement`, the second approach will automatically define the operators for `MyExtensionArray`.

A mixin class, `ExtensionScalarOpsMixin` supports this second approach. If developing an `ExtensionArray` subclass, for example `MyExtensionArray`, can simply include `ExtensionScalarOpsMixin` as a parent class of `MyExtensionArray`, and then call the methods `_add_arithmetic_ops()` and/or `_add_comparison_ops()` to hook the operators into your `MyExtensionArray` class, as follows:

```
from pandas.api.extensions import ExtensionArray, ExtensionScalarOpsMixin

class MyExtensionArray(ExtensionArray, ExtensionScalarOpsMixin):
    pass

MyExtensionArray._add_arithmetic_ops()
MyExtensionArray._add_comparison_ops()
```

Note: Since pandas automatically calls the underlying operator on each element one-by-one, this might not be as performant as implementing your own version of the associated operators directly on the `ExtensionArray`.

For arithmetic operations, this implementation will try to reconstruct a new `ExtensionArray` with the result of the element-wise operation. Whether or not that succeeds depends on whether the operation returns a result that's valid for the `ExtensionArray`. If an `ExtensionArray` cannot be reconstructed, an `ndarray` containing the scalars returned instead.

For ease of implementation and consistency with operations between pandas and NumPy `ndarrays`, we recommend *not* handling `Series` and `Indexes` in your binary ops. Instead, you should detect these cases and return `NotImplemented`. When pandas encounters an operation like `op(Series, ExtensionArray)`, pandas will

1. unbox the array from the `Series` (`Series.array`)
2. call `result = op(values, ExtensionArray)`
3. re-box the result in a `Series`

NumPy Universal Functions

Series implements `__array_ufunc__`. As part of the implementation, pandas unboxes the `ExtensionArray` from the *Series*, applies the ufunc, and re-boxes it if necessary.

If applicable, we highly recommend that you implement `__array_ufunc__` in your extension array to avoid coercion to an ndarray. See [the numpy documentation](#) for an example.

As part of your implementation, we require that you defer to pandas when a pandas container (*Series*, *DataFrame*, *Index*) is detected in `inputs`. If any of those is present, you should return `NotImplemented`. Pandas will take care of unboxing the array from the container and re-calling the ufunc with the unwrapped input.

Testing extension arrays

We provide a test suite for ensuring that your extension arrays satisfy the expected behavior. To use the test suite, you must provide several pytest fixtures and inherit from the base test class. The required fixtures are found in <https://github.com/pandas-dev/pandas/blob/master/pandas/tests/extension/conftest.py>.

To use a test, subclass it:

```
from pandas.tests.extension import base

class TestConstructors(base.BaseConstructorsTests):
    pass
```

See https://github.com/pandas-dev/pandas/blob/master/pandas/tests/extension/base/__init__.py for a list of all the tests available.

Compatibility with Apache Arrow

An `ExtensionArray` can support conversion to / from `pyarrow` arrays (and thus support for example serialization to the Parquet file format) by implementing two methods: `ExtensionArray.__arrow_array__` and `ExtensionDtype.__from_arrow__`.

The `ExtensionArray.__arrow_array__` ensures that `pyarrow` knows how to convert the specific extension array into a `pyarrow.Array` (also when included as a column in a pandas `DataFrame`):

```
class MyExtensionArray(ExtensionArray):
    ...

    def __arrow_array__(self, type=None):
        # convert the underlying array values to a pyarrow Array
        import pyarrow
        return pyarrow.array(..., type=type)
```

The `ExtensionDtype.__from_arrow__` method then controls the conversion back from `pyarrow` to a pandas `ExtensionArray`. This method receives a `pyarrow.Array` or `ChunkedArray` as only argument and is expected to return the appropriate pandas `ExtensionArray` for this dtype and the passed values:

```
class ExtensionDtype:
    ...

    def __from_arrow__(self, array: pyarrow.Array/ChunkedArray) -> ExtensionArray:
        ...
```

See more in the [Arrow documentation](#).

Those methods have been implemented for the nullable integer and string extension dtypes included in pandas, and ensure roundtrip to pyarrow and the Parquet file format.

4.5.3 Subclassing pandas data structures

Warning: There are some easier alternatives before considering subclassing pandas data structures.

1. Extensible method chains with *pipe*
2. Use *composition*. See [here](#).
3. Extending by *registering an accessor*
4. Extending by *extension type*

This section describes how to subclass pandas data structures to meet more specific needs. There are two points that need attention:

1. Override constructor properties.
2. Define original properties

Note: You can find a nice example in [geopandas](#) project.

Override constructor properties

Each data structure has several *constructor properties* for returning a new data structure as the result of an operation. By overriding these properties, you can retain subclasses through pandas data manipulations.

There are 3 constructor properties to be defined:

- `_constructor`: Used when a manipulation result has the same dimensions as the original.
- `_constructor_sliced`: Used when a manipulation result has one lower dimension(s) as the original, such as DataFrame single columns slicing.
- `_constructor_expanddim`: Used when a manipulation result has one higher dimension as the original, such as `Series.to_frame()`.

Following table shows how pandas data structures define constructor properties by default.

Property Attributes	Series	DataFrame
<code>_constructor</code>	Series	DataFrame
<code>_constructor_sliced</code>	NotImplementedError	Series
<code>_constructor_expanddim</code>	DataFrame	NotImplementedError

Below example shows how to define `SubclassedSeries` and `SubclassedDataFrame` overriding constructor properties.

```
class SubclassedSeries(pd.Series):  
  
    @property
```

(continues on next page)

(continued from previous page)

```

def _constructor(self):
    return SubclassedSeries

@property
def _constructor_expanddim(self):
    return SubclassedDataFrame

class SubclassedDataFrame(pd.DataFrame):

    @property
    def _constructor(self):
        return SubclassedDataFrame

    @property
    def _constructor_sliced(self):
        return SubclassedSeries

>>> s = SubclassedSeries([1, 2, 3])
>>> type(s)
<class '__main__.SubclassedSeries'>

>>> to_framed = s.to_frame()
>>> type(to_framed)
<class '__main__.SubclassedDataFrame'>

>>> df = SubclassedDataFrame({'A': [1, 2, 3], 'B': [4, 5, 6], 'C': [7, 8, 9]})
>>> df
   A  B  C
0  1  4  7
1  2  5  8
2  3  6  9

>>> type(df)
<class '__main__.SubclassedDataFrame'>

>>> sliced1 = df[['A', 'B']]
>>> sliced1
   A  B
0  1  4
1  2  5
2  3  6

>>> type(sliced1)
<class '__main__.SubclassedDataFrame'>

>>> sliced2 = df['A']
>>> sliced2
0    1
1    2
2    3
Name: A, dtype: int64

>>> type(sliced2)
<class '__main__.SubclassedSeries'>

```

Define original properties

To let original data structures have additional properties, you should let pandas know what properties are added. pandas maps unknown properties to data names overriding `__getattrute__`. Defining original properties can be done in one of 2 ways:

1. Define `_internal_names` and `_internal_names_set` for temporary properties which WILL NOT be passed to manipulation results.
2. Define `_metadata` for normal properties which will be passed to manipulation results.

Below is an example to define two original properties, “internal_cache” as a temporary property and “added_property” as a normal property

```
class SubclassedDataFrame2(pd.DataFrame):

    # temporary properties
    _internal_names = pd.DataFrame._internal_names + ['internal_cache']
    _internal_names_set = set(_internal_names)

    # normal properties
    _metadata = ['added_property']

    @property
    def _constructor(self):
        return SubclassedDataFrame2

>>> df = SubclassedDataFrame2({'A': [1, 2, 3], 'B': [4, 5, 6], 'C': [7, 8, 9]})
>>> df
   A  B  C
0  1  4  7
1  2  5  8
2  3  6  9

>>> df.internal_cache = 'cached'
>>> df.added_property = 'property'

>>> df.internal_cache
cached
>>> df.added_property
property

# properties defined in _internal_names is reset after manipulation
>>> df[['A', 'B']].internal_cache
AttributeError: 'SubclassedDataFrame2' object has no attribute 'internal_cache'

# properties defined in _metadata are retained
>>> df[['A', 'B']].added_property
property
```

4.5.4 Plotting backends

Starting in 0.25 pandas can be extended with third-party plotting backends. The main idea is letting users select a plotting backend different than the provided one based on Matplotlib. For example:

```
>>> pd.set_option('plotting.backend', 'backend.module')
>>> pd.Series([1, 2, 3]).plot()
```

This would be more or less equivalent to:

```
>>> import backend.module
>>> backend.module.plot(pd.Series([1, 2, 3]))
```

The backend module can then use other visualization tools (Bokeh, Altair, ...) to generate the plots.

Libraries implementing the plotting backend should use `entry points` to make their backend discoverable to pandas. The key is "pandas_plotting_backends". For example, pandas registers the default "matplotlib" backend as follows.

```
# in setup.py
setup( # noqa: F821
    ...,
    entry_points={
        "pandas_plotting_backends": [
            "matplotlib = pandas:plotting._matplotlib",
        ],
    },
)
```

More information on how to implement a third-party plotting backend can be found at https://github.com/pandas-dev/pandas/blob/master/pandas/plotting/_init__.py#L1.

4.6 Developer

This section will focus on downstream applications of pandas.

4.6.1 Storing pandas DataFrame objects in Apache Parquet format

The `Apache Parquet` format provides key-value metadata at the file and column level, stored in the footer of the Parquet file:

```
5: optional list<KeyValue> key_value_metadata
```

where `KeyValue` is

```
struct KeyValue {
  1: required string key
  2: optional string value
}
```

So that a `pandas.DataFrame` can be faithfully reconstructed, we store a `pandas` metadata key in the `FileMetaData` with the value stored as :


```
{'index_columns': [<descr0>, <descr1>, ...],
 'column_indexes': [<ci0>, <ci1>, ..., <ciN>],
 'columns': [<c0>, <c1>, ...],
 'pandas_version': $VERSION,
 'creator': {
   'library': $LIBRARY,
   'version': $LIBRARY_VERSION
 }}
```

The “descriptor” values <descr0> in the 'index_columns' field are strings (referring to a column) or dictionaries with values as described below.

The <c0>/<ci0> and so forth are dictionaries containing the metadata for each column, *including the index columns*. This has JSON form:

```
{'name': column_name,
 'field_name': parquet_column_name,
 'pandas_type': pandas_type,
 'numpy_type': numpy_type,
 'metadata': metadata}
```

See below for the detailed specification for these.

Index Metadata Descriptors

RangeIndex can be stored as metadata only, not requiring serialization. The descriptor format for these as is follows:

```
index = pd.RangeIndex(0, 10, 2)
{'kind': 'range',
 'name': index.name,
 'start': index.start,
 'stop': index.stop,
 'step': index.step}
```

Other index types must be serialized as data columns along with the other DataFrame columns. The metadata for these is a string indicating the name of the field in the data columns, for example '__index_level_0__'.

If an index has a non-None name attribute, and there is no other column with a name matching that value, then the index.name value can be used as the descriptor. Otherwise (for unnamed indexes and ones with names colliding with other column names) a disambiguating name with pattern matching __index_level_\d+__ should be used. In cases of named indexes as data columns, name attribute is always stored in the column descriptors as above.

Column Metadata

pandas_type is the logical type of the column, and is one of:

- Boolean: 'bool'
- Integers: 'int8', 'int16', 'int32', 'int64', 'uint8', 'uint16', 'uint32', 'uint64'
- Floats: 'float16', 'float32', 'float64'
- Date and Time Types: 'datetime', 'datetime64', 'timedelta'
- String: 'unicode', 'bytes'
- Categorical: 'categorical'