

Pandas

Goals

Fund the continued maintenance and improvements for Pandas (McKinney (2010)), a Python library providing high-performance, easy-to-use data structures. Pandas is a foundational library in the Scientific Python Ecosystem, providing the most widely-used data structure for heterogenous, tabular data ().

Library Maintenance

Pandas is a large library with many users. We have many open issues (about 3,000) and Pull Requests (about 100). Keeping up with the stream of updates is time consuming, and as such we would like to fund time dedicated specifically for maintenance. This would include

1. Promptly triaging newly opened issues
2. Promptly reviewing new pull requests
3. Ensuring issue conversations are progressing
4. Ensuring that Pull Requests are not going stale and being abandoned
5. Triage newly opened issues. Bug reports need to be verified and assigned a priority. Feature requests need to be put through an initial design discussion, ensuring that the relevant community members have a chance to guide the design.
6. Reviewing Pull Requests. Newly opened contributions should be given a prompt, high-level overview. Relevant experts in that code-section should be asked to provide their review.
7. Follow up on issues and pull requests that are going stale. Pandas enjoys a large number of one-time pull requests from a member of the community. These contributors are often unable to respond quickly to multiple rounds of review. This can result in the contribution going stale and the Pull Request abandoned.

Native String Data Type

Currently, pandas stores text data in an `object`-dtype NumPy array. Each array stores Python strings. While pragmatic, since we rely on NumPy for storage and Python for string operations, this is memory inefficient and slow. We'd like to provide a native string type for pandas.

The most obvious alternative is Apache Arrow. Arrow's data model is well-suited to storing arrays of textual data. Arrow arrays can already be stored inside a DataFrame using the extension array interface. We propose the following

1. Work with the Arrow developers to implement the string operations on Arrow memory. This may be done in the Arrow C++ library (benefiting all libraries using it, including Python, R, and Ruby), or in a Python-specific manor using Numba (Lam, Pitrou, and Seibert (2015)) (following the Fletcher proof of concept).
2. Provide a backwards-compatible transition for current pandas users. This large of a change (including a new dependency) may require a period where users must opt into the new behavior.

Documentation Validation

To improve the quality and consistency of Pandas documentation, we've developed tooling to check docstrings in a variety of ways. Every docstring is checked for

1. Consistency: Ensuring that every docstring has the same components (parameters, return values, notes, examples, etc.). Ensuring that each section is formatted correctly.
2. Completeness: Ensuring that every function parameter is documented.
3. Correctness: Ensuring that the examples run correctly.

When a user submits a Pull Request to pandas, their changes to the documentation are automatically checked, and informative error messages inform them of any issues.

Like many other projects, pandas uses the `numpydoc` standard for writing docstrings. With the collaboration of the `numpydoc` maintainers, we'd like to move the project-agnostic tooling we've written to a different project that any `numpydoc`-using project can benefit from.

If possible, we'd like this project to be undertaken by a member of an unrepresented minority, with mentorship provided by the pandas maintainers. This project primarily requires experience with *using* pandas, NumPy, and related libraries, rather than deep knowledge of pandas' internals. Pandas has a diverse community, just not at the maintainer level (yet).

Performance Monitoring

Pandas uses airspeed velocity to monitor for performance regressions. ASV itself is a fabulous tool, but requires some additional work to be integrated into an open source project’s workflow.

The asv-runner organization, currently made up of pandas maintainers, provides tools built on top of ASV. We have a physical machine for running a number of project’s benchmarks and tools managing the benchmark runs and reporting on results.

We’d like to fund improvements and maintenance of these tools to

- Be more reliable. Currently, they’re maintained on the nights and weekends when a maintainer has free time. This occasionally results in days or weeks of downtime.
- Tune the system to improve benchmark stability
- Report performance regressions to a project rather than relying on maintainers manually checking for regressions.
- Build a GitHub bot to request ASV runs *before* a Pull Request is merged. The benchmarks are too expensive to run as part of every commit. Running on-demand from a maintainer provides a nice balance.
- Pay for hosting a dedicated benchmark server. A maintainer’s basement isn’t the most stable environment for a machine.

This project has an impact beyond pandas. We run the benchmarks for many foundational libraries in the Scientific Python Ecosystem including xarray, Dask, scikit-image, scikit-learn, and PyMC3. Each of these projects would benefit from improvements made to the tooling.

Existing Support

Pandas currently has about two FTE.

- Anaconda Inc funds Tom Augspurger (50% FTE) and Brock Mendel (100% FTE)
- Ursa Labs funds Joris van den Bossche (50% FTE).
- NumFocus Small Development Grant for “Improving and modernizing the introductory *Getting Started* pages of the pandas documentation” (\$5,000)

An up-to-date list of current members and institutional partners can be found in our governance documents.

References

Lam, Siu Kwan, Antoine Pitrou, and Stanley Seibert. 2015. “Numba: A Llvm-Based Python Jit Compiler.” In *Proceedings of the Second Workshop on the*

Llvm Compiler Infrastructure in Hpc, 7. ACM.

McKinney, Wes. 2010. “Data Structures for Statistical Computing in Python.”
In *Proceedings of the 9th Python in Science Conference*, edited by Stéfan van
der Walt and Jarrod Millman, 51–56.