
xarray*extras* Documentation

Release 0.2.1

xarray*extras* Developers

2018-07-22

Contents

1 Installation	1
1.1 Required dependencies	1
1.2 Testing	1
2 What's New	3
2.1 v0.2.1 (2018-07-22)	3
2.2 v0.2.0 (2018-07-15)	3
2.3 v0.1.0 (2018-05-19)	3
3 API Reference	5
3.1 csv	5
3.2 cumulatives	6
3.3 interpolate	7
3.4 numba_extras	8
3.5 sort	9
4 License	11
Python Module Index	13

CHAPTER 1

Installation

1.1 Required dependencies

- Python 3.5 or later
- [scipy](#)
- [xarray](#)
- [dask](#)
- [numba](#)
- C compiler (only if building from sources)

1.2 Testing

To run the test suite after installing `xarray_extras`, first install (via pypi or conda)

- [py.test](#): Simple unit testing library

and run `py.test --pyargs xarray_extras`.

CHAPTER 2

What's New

2.1 v0.2.1 (2018-07-22)

- Added parameter nogil=True to `xarray_extras.csv.to_csv()`, which will switch to a C-accelerated implementation instead of pandas `to_csv` (albeit with caveats). Fixed deadlock in `to_csv` as well as compatibility with dask distributed. Pandas code (when using `nogil=False`) is not wrapped by a subprocess anymore, which means it won't be able to use more than 1 CPU (but compression can run in pipeline). `to_csv` has lost the ability to write to a buffer - only file paths are supported now.
- AppVeyor integration

2.2 v0.2.0 (2018-07-15)

- New function `xarray_extras.csv.to_csv()`
- Speed up interpolation for k=2 and k=3
- CI: Rigorous tracking of minimum dependency versions
- CI: Explicit support for Python 3.7

2.3 v0.1.0 (2018-05-19)

Initial release.

CHAPTER 3

API Reference

3.1 csv

CSV file type support

```
xarray_extras.csv.to_csv(x, path, *, nogil=True, **kwargs)  
    Print DataArray to CSV.
```

When x has numpy backend, this function is equivalent to:

```
x.to_pandas().to_csv(path_or_buf, **kwargs)
```

When x has dask backend, this function returns a dask delayed object which will write to the disk only when its .compute() method is invoked.

Formatting and optional compression are parallelised across all available CPUs, using one dask task per chunk on the first dimension. Chunks on other dimensions will be merged ahead of computation.

Parameters

- **x** – xarray.DataArray with one or two dimensions
- **path (str)** – Output file path
- **nogil (bool)** – If True, use accelerated C implementation. Several kwargs won't be processed correctly (see limitations below). If False, use pandas to_csv method (slow, and does not release the GIL). nogil=True exclusively supports float and integer values dtypes (but the coords can be anything). In case of incompatible dtype, nogil is automatically switched to False.
- **kwargs** – Passed verbatim to `pandas.DataFrame.to_csv()` or `pandas.Series.to_csv()`

Limitations

- Fancy URIs are not (yet) supported.
- compression='zip' is not supported. All other compression methods (gzip, bz2, xz) are supported.

- When running with nogil=True, the following parameters are ignored: columns, quoting, quotechar, doublequote, escapechar, chunksize, decimal

Distributed

This function supports `dask distributed`, with the caveat that all workers must write to the same shared mount-point and that the shared filesystem must strictly guarantee **close-open coherency**, meaning that one must be able to call `write()` and then `close()` on a file descriptor from one host and then immediately afterwards `open()` from another host and see the output from the first host. Note that, for performance reasons, most network filesystems do not enable this feature by default.

Alternatively, one may write to local mountpoints and then manually collect and concatenate the partial outputs.

3.2 cumulatives

xarray extensions to cumsum and cumprod

`xarray_extras.cumulatives.cummean(x, dim, skipna=None)`

$$y_i = \text{mean}(x_0, x_1, \dots x_i)$$

Parameters

- `x` – any xarray object
- `dim (str)` – dimension along which to calculate the mean
- `skipna (bool)` – If True, skip missing values (as marked by NaN). By default, only skips missing values for float dtypes; other dtypes either do not have a sentinel missing value (int) or skipna=True has not been implemented (object, datetime64 or timedelta64).

`xarray_extras.cumulatives.compound_sum(x, c, xdim, cdim)`

Compound sum on arbitrary points of x along dim.

Parameters

- `x` – Any xarray object containing the data to be compounded
- `c (xarray.DataArray)` – array where every row contains elements of `x.coords[xdim]` and is used to build a point of the output. The cells in the row are matched against `x.coords[dim]` and perform a sum. If different rows of `c` require different amounts of points from `x`, they must be padded on the right with NaN, NaT, or “” (respectively for numbers, datetimes, and strings).
- `xdim (str)` – dimension of `x` to acquire data from. The coord associated to it must be monotonic ascending.
- `cdim (str)` – dimension of `c` that represent the vector of points to be compounded for every point of `dim`

Returns DataArray with all dims from `x` and `c`, except `xdim` and `cdim`, and the same dtype as `x`.

example:

```
>>> x = xarray.DataArray(  
>>>     [10, 20, 30],  
>>>     dims=['x'], coords={'x': ['foo', 'bar', 'baz']})  
>>> c = xarray.DataArray(  
>>>     [['foo', 'baz', None],  
>>>      ['bar', 'baz', 'baz']],
```

(continues on next page)

(continued from previous page)

```
>>> dims=['y', 'c'], coords={'y': ['new1', 'new2']})
>>> compound_sum(x, c, 'x', 'c')
<xarray.DataArray (y: 2)>
array([40, 80])
Coordinates:
* y          (y) <U4 'new1' 'new2'
```

`xarray_extras.cumulatives.compound_prod(x, c, xdim, cdim)`
 Compound product among arbitrary points of x along dim See [compound_sum\(\)](#).

`xarray_extras.cumulatives.compound_mean(x, c, xdim, cdim)`
 Compound mean among arbitrary points of x along dim See [compound_sum\(\)](#).

3.3 interpolate

xarray interpolation functions

`xarray_extras.interpolate.splrep(a, dim, k=3)`
 Calculate the univariate B-spline for an N-dimensional array

Parameters

- `a (xarray.DataArray)` – any `DataArray`
- `dim` – dimension of a to be interpolated. `a.coords[dim]` must be strictly monotonic ascending. All int, float (not complex), or datetime dtypes are supported.
- `k (int)` – B-spline order:

k	interpolation kind
0	nearest neighbour
1	linear
2	quadratic
3	cubic

Returns `Dataset` with t, c, k (knots, coefficients, order) variables, the same shape and coords as the input, that can be passed to `splev()`.

Example:

```
>>> x = np.arange(0, 120, 20)
>>> x = xarray.DataArray(x, dims=['x'], coords={'x': x})
>>> s = xarray.DataArray(np.linspace(1, 20, 5), dims=['s'])
>>> y = np.exp(-x / s)
>>> x_new = np.arange(0, 120, 1)
>>> tck = splrep(y, 'x')
>>> y_new = splev(x_new, tck)
```

Features

- Interpolate a ND array on any arbitrary dimension
- dask supported on both on the interpolated array and `x_new`
- Supports ND `x_new` arrays

- The CPU-heavy interpolator generation (`splrep()`) is executed only once and then can be applied to multiple `x_new` (`splev()`)
- memory-efficient
- Can be pickled and used on dask distributed

Limitations

- Chunks are not supported along dim on the interpolated dimension.

`xarray_extras.interpolate.splev(x_new, tck, extrapolate=True)`
Evaluate the B-spline generated with `splrep()`.

Parameters

- **`x_new`** – Any `DataArray` with any number of dims, not necessarily the original interpolation dim. Alternatively, it can be any 1-dimensional array-like; it will be automatically converted to a `DataArray` on the interpolation dim.

- **`tck` (`xarray.Dataset`)** – As returned by `splrep()`. It can have been:
 - transposed (not recommended, as performance will drop if c is not C-contiguous)
 - sliced, reordered, or (re)chunked, on any dim except the interpolation dim
 - computed from dask to numpy backend
 - round-tripped to disk

- **`extrapolate` –**

True Extrapolate the first and last polynomial pieces of b-spline functions active on the base interval

False Return NaNs outside of the base interval

'periodic' Periodic extrapolation is used

'clip' Return `y[0]` and `y[-1]` outside of the base interval

Returns `DataArray` with all dims of the interpolated array, minus the interpolation dim, plus all dims of `x_new`

See `splrep()` for usage example.

3.4 numba_extras

Extensions to numba

`xarray_extras.numba_extras.guvectorize(signature, layout, **kwds)`

Convenience wrapper around `numba.guvectorize()`. Generate signature for all possible data types and set a few healthy defaults.

Parameters

- **`signature` (`str`)** – numba signature, containing {T}
- **`layout` (`str`)** – as in `numba.guvectorize()`
- **`kwds`** – passed verbatim to `numba.guvectorize()`. This function changes the default for cache from False to True.

example:

```
guvectorize("{T}[:, {T}[:, {(i)->(i)}"]
```

Is the same as:

```
numba.guvectorize([
    "float32[:, float32[:]]",
    "float64[:, float64[:]]",
    ...
], "(i)->(i)", cache=True)
```

Note: Discussing upstream fix; see <https://github.com/numba/numba/issues/2936>.

3.5 sort

Sorting functions

`xarray_extras.sort.topk(a, k, dim, split_every=None)`

Extract the k largest elements from a on the given dimension, and return them sorted from largest to smallest. If k is negative, extract the -k smallest elements instead, and return them sorted from smallest to largest.

This assumes that k is small. All results will be returned in a single chunk along the given axis.

`xarray_extras.sort.argsort(a, k, dim, split_every=None)`

Extract the indexes of the k largest elements from a on the given dimension, and return them sorted from largest to smallest. If k is negative, extract the -k smallest elements instead, and return them sorted from smallest to largest.

This assumes that k is small. All results will be returned in a single chunk along the given axis.

`xarray_extras.sort.take_along_dim(a, ind, dim)`

Use the output of `argsort()` to pick points from a.

Parameters

- `a` – any xarray object
- `ind` – array of ints, as returned by `argsort()`
- `dim` – dimension along which `argsort` was executed

An example that uses all of the above functions is *source attribution*. Given a generic function $y = f(x_0, x_1, \dots, x_i)$, which is embarrassingly parallel along a given dimension, one wants to find:

- the top k elements of y along the dimension
- the elements of all x's that generated the top k elements of y

```
>>> from xarray import DataArray
>>> from xarray_extras.sort import *
>>> x = DataArray([[5, 3, 2, 8, 1],
>>>                 [0, 7, 1, 3, 2]], dims=['x', 's'])
>>> y = x.sum('x') # y = f(x), embarrassingly parallel among dimension 's'
>>> y
<xarray.DataArray (s: 5)>
array([ 5, 10,  3, 11,  3])
Dimensions without coordinates: s
>>> top_y = topk(y, 3, 's')
```

(continues on next page)

(continued from previous page)

```
>>> top_y
<xarray.DataArray (s: 3)>
array([11, 10,  5])
Dimensions without coordinates: s
>>> top_x = take_along_dim(x, argtopk(y, 3, 's'), 's')
>>> top_x
<xarray.DataArray (x: 2, s: 3)>
array([[8, 3, 5],
       [3, 7, 0]])
Dimensions without coordinates: x, s
```

CHAPTER 4

License

xarray_extras is available under the open source [LGPL License](#).

Python Module Index

X

`xarray_extras.csv`, [5](#)
`xarray_extras.cumulatives`, [6](#)
`xarray_extras.interpolate`, [7](#)
`xarray_extras.numba_extras`, [8](#)
`xarray_extras.sort`, [9](#)

A

`argtopk()` (in module `xarray_extras.sort`), 9

C

`compound_mean()` (in module `xarray_extras.cumulatives`), 7
`compound_prod()` (in module `xarray_extras.cumulatives`),
7
`compound_sum()` (in module `xarray_extras.cumulatives`),
6
`cummean()` (in module `xarray_extras.cumulatives`), 6

G

`guvectorize()` (in module `xarray_extras.numba_extras`), 8

S

`splev()` (in module `xarray_extras.interpolate`), 8
`splrep()` (in module `xarray_extras.interpolate`), 7

T

`take_along_dim()` (in module `xarray_extras.sort`), 9
`to_csv()` (in module `xarray_extras.csv`), 5
`topk()` (in module `xarray_extras.sort`), 9

X

`xarray_extras.csv` (module), 5
`xarray_extras.cumulatives` (module), 6
`xarray_extras.interpolate` (module), 7
`xarray_extras.numba_extras` (module), 8
`xarray_extras.sort` (module), 9