
clisops Documentation

Release 0.9.3

Elle Smith

Sep 30, 2022

CONTENTS:

1 Quick Guide	1
1.1 clisops - climate simulation operations	1
2 Installation	3
2.1 Stable release	3
2.2 From sources	3
3 Usage	5
4 API	7
4.1 Core subset functionality	7
4.2 Core average functionality	16
4.3 Subset operation	18
4.4 Average operations	19
4.5 Common functions	21
4.6 Output Utilities	21
4.7 File Namers	22
4.8 Dataset Utilities	22
4.9 Time Utilities	24
5 Examples	25
5.1 Subsetting	25
5.2 Subsetting utilities	32
5.3 Averaging over dimensions of the dataset	42
6 Contributing	51
6.1 Types of Contributions	51
6.2 Get Started!	52
6.3 Pull Request Guidelines	53
6.4 Tips	54
6.5 Versioning	54
6.6 Deploying	54
6.7 Packaging	54
7 Credits	55
7.1 Development Lead	55
7.2 Co-Developers	55
7.3 Contributors	55
8 Version History	57
8.1 v0.9.3 (2022-10-03)	57

8.2	Bug Fixes	57
8.3	Other Changes	57
8.4	v0.9.2 (2022-09-06)	57
8.5	v0.9.1 (2022-05-12)	58
8.6	v0.9.0 (2022-04-13)	58
8.7	v0.8.0 (2022-01-13)	59
8.8	v0.7.0 (2021-10-26)	60
8.9	v0.6.5 (2021-06-10)	60
8.10	v0.6.4 (2021-05-17)	61
8.11	v0.6.3 (2021-03-30)	61
8.12	v0.6.2 (2021-03-22)	62
8.13	v0.6.1 (2021-02-23)	62
8.14	v0.6.0 (2021-02-22)	62
8.15	v0.5.1 (2021-01-11)	63
8.16	v0.5.0 (2020-12-17)	63
8.17	v0.4.0 (2020-11-10)	64
8.18	v0.3.1 (2020-08-04)	65
8.19	v0.3.0 (2020-07-23)	65
8.20	v0.2.1 (2020-07-08)	65
8.21	v0.2.0 (2020-06-19)	65
8.22	v0.1.0 (2020-04-22)	66
9	Indices and tables	67
	Python Module Index	69
	Index	71

1.1 clisops - climate simulation operations

The `clisops` package (pronounced “clie-sops”) provides a python library for running *data-reduction* operations on `Xarray` data sets or files that can be interpreted by `Xarray`. These basic operations (subsetting, averaging and regridding) are likely to work where data structures are NetCDF-centric, such as those found in ESGF data sets.

`clisops` is employed by the `daops` library to perform its basic operations once `daops` has applied any necessary *fixes* to data in order to remove irregularities/anomalies. Users are recommended to investigate using `daops` directly in order to access these *fixes* which may affect the scientific credibility of the results.

`clisops` can be used stand-alone to read individual, or groups of, NetCDF files directly.

- Free software: BSD
- Documentation: <https://clisops.readthedocs.io>.

1.1.1 Features

The package provides the following operations:

- subset
- average
- regrid

1.1.2 Credits

This package was created with `Cookiecutter` and the `audreyr/cookiecutter-pypackage` project template.

- `Cookiecutter`: <https://github.com/audreyr/cookiecutter>
- `cookiecutter-pypackage`: <https://github.com/audreyr/cookiecutter-pypackage>

INSTALLATION

2.1 Stable release

To install `clisops`, run this command in your terminal:

```
$ pip install clisops
```

This is the preferred method to install `clisops`, as it will always install the most recent stable release.

If you don't have `pip` installed, this [Python installation guide](#) can guide you through the process.

Warning: Some average operations (`clisops.core.average_shape`) require at least version 0.5.2 of the `xESMF` package. Unfortunately, this package is not available on pypi at the time these lines were written and it also depends on packages (`ESMF`, `ESMpy`) unavailable on windows. It can still be installed on osx/linux through `conda` or directly [from source](<https://github.com/pangeo-data/xESMF/>).

Another optional dependency is `[pygeos]`(<https://pygeos.readthedocs.io/en/latest/index.html>). If installed, the performance of `core.subset.create_mask` and `core.subset.subset_shape` are greatly improved.

2.2 From sources

The sources for `clisops` can be downloaded from the [Github repo](#).

You can either clone the public repository:

```
$ git clone git://github.com/roocs/clisops
```

Create Conda environment named `clisops`.

Warning: For windows users, you might need to edit the `environment.yml` file and remove the `xESMF` package.

```
$ conda env create -f environment.yml  
$ source activate clisops
```

Install `clisops` in development mode:

```
$ pip install -r requirements.txt
$ pip install -r requirements_dev.txt
$ python setup.py develop
```

Run tests:

```
$ pytest -v tests/
```


USAGE

To use clisops in a project:

```
import clisops
```

For information on the configuration options available in clisops, see: <https://roocs-utils.readthedocs.io/en/latest/configuration.html#clisops>

4.1 Core subset functionality

Subset module.

```
clisops.core.subset.create_mask(*, x_dim: DataArray, y_dim: DataArray, poly: GeoDataFrame,  
                                wrap_lons: bool = False, check_overlap: bool = False)
```

Create a mask with values corresponding to the features in a GeoDataFrame using vectorize methods.

The returned mask's points have the value of the first geometry of *poly* they fall in.

Parameters

- **x_dim** (*xarray.DataArray*) – X or longitudinal dimension of xarray object. Can also be given through *ds_in*.
- **y_dim** (*xarray.DataArray*) – Y or latitudinal dimension of xarray object. Can also be given through *ds_in*.
- **poly** (*gpd.GeoDataFrame*) – A GeoDataFrame used to create the xarray.DataArray mask. If its index doesn't have an integer dtype, it will be reset to integers, which will be used in the mask.
- **wrap_lons** (*bool*) – Shift vector longitudes by -180,180 degrees to 0,360 degrees; Default = False
- **check_overlap** (*bool*) – Perform a check to verify if shapes contain overlapping geometries.

Returns

xarray.DataArray

Examples

```
import geopandas as gpd  
from clisops.core.subset import create_mask  
ds = xr.open_dataset(path_to_tasmin_file)  
polys = gpd.read_file(path_to_multi_shape_file)  
  
# Get a mask from all polygons in the shape file  
mask = create_mask(x_dim=ds.lon, y_dim=ds.lat, poly=polys)  
ds = ds.assign_coords(regions=mask)  
  
# Operations can be applied to each region with `groupby`. Ex:
```

(continues on next page)

(continued from previous page)

```

ds = ds.groupby('regions').mean()

# Extra step to retrieve the names of those polygons stored in another column (here
# → "id")
region_names = xr.DataArray(polys.id, dims=('regions',))
ds = ds.assign_coords(regions_names=region_names)

```

`clisops.core.subset.create_weight_masks(ds_in: Union[DataArray, Dataset], poly: GeoDataFrame)`

Create weight masks corresponding to the features in a GeoDataFrame using xESMF.

The returned masks values are the fraction of the corresponding polygon's area that is covered by the grid cell. Summing along the spatial dimension will give 1 for each geometry. Requires xESMF >= 0.5.0.

Parameters

- **ds_in** (*Union[xarray.DataArray, xarray.Dataset]*) – xarray object containing the grid information, as understood by xESMF. For 2D lat/lon coordinates, the bounds arrays are required,
- **poly** (*gpd.GeoDataFrame*) – GeoDataFrame used to create the xarray.DataArray mask. One mask will be created for each row in the dataframe. Will be converted to EPSG:4326 if needed.

Returns

xarray.DataArray – Has a new *geom* dimension corresponding to the index of the input GeoDataframe. Non-geometry columns of *poly* are copied as auxiliary coordinates.

Examples

```

>>> import geopandas as gpd
>>> import xarray as xr
>>> from clisops.core.subset import create_weight_masks
>>> ds = xr.open_dataset(path_to_tasmin_file)
>>> polys = gpd.read_file(path_to_multi_shape_file)
...
# Get a weight mask for each polygon in the shape file
>>> mask = create_weight_masks(x_dim=ds.lon, y_dim=ds.lat, poly=polys)

```

`clisops.core.subset.distance(da: Union[DataArray, Dataset], *, lon: Union[float, Sequence[float], DataArray], lat: Union[float, Sequence[float], DataArray])`

Return distance to a point in meters.

Parameters

- **da** (*Union[xarray.DataArray, xarray.Dataset]*) – Input data.
- **lon** (*Union[float, Sequence[float], xarray.DataArray]*) – Longitude coordinate.
- **lat** (*Union[float, Sequence[float], xarray.DataArray]*) – Latitude coordinate.

Returns

xarray.DataArray – Distance in meters to point.

Examples

```
import xarray as xr
from clisops.core.subset import distance

# To get the indices from the closest point, use:
da = xr.open_dataset(path_to_pr_file).pr
d = distance(da, lon=-75, lat=45)
k = d.argmin()
i, j, _ = np.unravel_index(k, d.shape)
```

`clisops.core.subset.shape_bbox_indexer(ds, poly)`

Return a spatial indexer that selects the indices of the grid cells covering the given geometries.

Parameters

- **ds** (*xr.Dataset*) – Input dataset.
- **poly** (*gpd.GeoDataFrame*) – Shapes to cover.

Returns

dict – xarray indexer along native dataset coordinates, to be used as an argument to *isel*.

Examples

```
>>> indexer = shape_bbox_indexer(ds, poly)
>>> ds.isel(indexer)
```

Notes

This is used in particular to restrict the domain of a dataset before computing the weights for a spatial average.

`clisops.core.subset.subset_bbox(da: Union[DataArray, Dataset], lon_bnds: Union[array, Tuple[Optional[float], Optional[float]]] = None, lat_bnds: Union[array, Tuple[Optional[float], Optional[float]]] = None, start_date: Optional[str] = None, end_date: Optional[str] = None, first_level: Optional[Union[int, float]] = None, last_level: Optional[Union[int, float]] = None, time_values: Optional[Sequence[str]] = None, level_values: Optional[Union[Sequence[float], Sequence[int]]] = None) → Union[DataArray, Dataset]`

Subset a DataArray or Dataset spatially (and temporally) using a lat lon bounding box and date selection.

Return a subset of a DataArray or Dataset for grid points falling within a spatial bounding box defined by longitude and latitudinal bounds and for dates falling within provided bounds.

TODO: returns the what? In the case of a lat-lon rectilinear grid, this simply returns the

Parameters

- **da** (*Union[xarray.DataArray, xarray.Dataset]*) – Input data.
- **lon_bnds** (*Union[np.array, Tuple[Optional[float], Optional[float]]]*) – List of minimum and maximum longitudinal bounds. Optional. Defaults to all longitudes in original data-array.
- **lat_bnds** (*Union[np.array, Tuple[Optional[float], Optional[float]]]*) – List of minimum and maximum latitudinal bounds. Optional. Defaults to all latitudes in original data-array.

- **start_date** (*Optional[str]*) – Start date of the subset. Date string format – can be year (“%Y”), year-month (“%Y-%m”) or year-month-day (“%Y-%m-%d”). Defaults to first day of input data-array.
- **end_date** (*Optional[str]*) – End date of the subset. Date string format – can be year (“%Y”), year-month (“%Y-%m”) or year-month-day (“%Y-%m-%d”). Defaults to last day of input data-array.
- **first_level** (*Optional[Union[int, float]]*) – First level of the subset. Can be either an integer or float. Defaults to first level of input data-array.
- **last_level** (*Optional[Union[int, float]]*) – Last level of the subset. Can be either an integer or float. Defaults to last level of input data-array.
- **time_values** (*Optional[Sequence[str]]*) – A list of datetime strings to subset.
- **level_values** (*Optional[Union[Sequence[float], Sequence[int]]]*) – A list of level values to select.

Returns

Union[xarray.DataArray, xarray.Dataset] – Subsetted *xarray.DataArray* or *xarray.Dataset*

Notes

subset_bbox expects the lower and upper bounds to be provided in ascending order. If the actual coordinate values are descending then this will be detected and your selection reversed before the data subset is returned.

Examples

```
import xarray as xr
from clisops.core.subset import subset_bbox

ds = xr.open_dataset(path_to_pr_file)

# Subset lat lon
prSub = subset_bbox(ds.pr, lon_bnds=[-75, -70], lat_bnds=[40, 45])
```

`clisops.core.subset.subset_gridpoint`(*da: Union[DataArray, Dataset], lon: Optional[Union[float, Sequence[float], DataArray]] = None, lat: Optional[Union[float, Sequence[float], DataArray]] = None, start_date: Optional[str] = None, end_date: Optional[str] = None, first_level: Optional[Union[int, float]] = None, last_level: Optional[Union[int, float]] = None, tolerance: Optional[float] = None, add_distance: bool = False*) → *Union[DataArray, Dataset]*

Extract one or more nearest gridpoint(s) from datarray based on lat lon coordinate(s).

Return a subsetted data array (or Dataset) for the grid point(s) falling nearest the input longitude and latitude coordinates. Optionally subset the data array for years falling within provided date bounds. Time series can optionally be subsetted by dates. If 1D sequences of coordinates are given, the gridpoints will be concatenated along the new dimension “site”.

Parameters

- **da** (*Union[xarray.DataArray, xarray.Dataset]*) – Input data.
- **lon** (*Optional[Union[float, Sequence[float], xarray.DataArray]]*) – Longitude coordinate(s). Must be of the same length as lat.

- **lat** (*Optional[Union[float, Sequence[float], xarray.DataArray]*) – Latitude coordinate(s). Must be of the same length as lon.
- **start_date** (*Optional[str]*) – Start date of the subset. Date string format – can be year (“%Y”), year-month (“%Y-%m”) or year-month-day (“%Y-%m-%d”). Defaults to first day of input data-array.
- **end_date** (*Optional[str]*) – End date of the subset. Date string format – can be year (“%Y”), year-month (“%Y-%m”) or year-month-day (“%Y-%m-%d”). Defaults to last day of input data-array.
- **first_level** (*Optional[Union[int, float]]*) – First level of the subset. Can be either an integer or float. Defaults to first level of input data-array.
- **last_level** (*Optional[Union[int, float]]*) – Last level of the subset. Can be either an integer or float. Defaults to last level of input data-array.
- **tolerance** (*Optional[float]*) – Masks values if the distance to the nearest gridpoint is larger than tolerance in meters.
- **add_distance** (*bool*)

Returns

Union[xarray.DataArray, xarray.Dataset] – Subsetted *xarray.DataArray* or *xarray.Dataset*

Examples

```
import xarray as xr
from clisops.core.subset import subset_gridpoint

ds = xr.open_dataset(path_to_pr_file)

# Subset lat lon point
prSub = subset_gridpoint(ds.pr, lon=-75, lat=45)

# Subset multiple variables in a single dataset
ds = xr.open_mfdataset([path_to_tasmax_file, path_to_tasmin_file])
dsSub = subset_gridpoint(ds, lon=-75, lat=45)
```

`clisops.core.subset.subset_level`(*da: Union[DataArray, Dataset], first_level: Optional[Union[int, float, str]] = None, last_level: Optional[Union[int, float, str]] = None*) → *Union[DataArray, Dataset]*

Subset input *DataArray* or *Dataset* based on first and last levels.

Return a subset of a *DataArray* or *Dataset* for levels falling within the provided bounds.

Parameters

- **da** (*Union[xarray.DataArray, xarray.Dataset]*) – Input data.
- **first_level** (*Optional[Union[int, float, str]]*) – First level of the subset (specified as the value, not the index). Can be either an integer or float. Defaults to first level of input data-array.
- **last_level** (*Optional[Union[int, float, str]]*) – Last level of the subset (specified as the value, not the index). Can be either an integer or float. Defaults to last level of input data-array.

Returns

Union[xarray.DataArray, xarray.Dataset] – Subsetted *xarray.DataArray* or *xarray.Dataset*

Examples

```
import xarray as xr
from clisops.core.subset import subset_level

ds = xr.open_dataset(path_to_pr_file)

# Subset complete levels
prSub = subset_level(ds.pr, first_level=0, last_level=30)

# Subset single level
prSub = subset_level(ds.pr, first_level=1000, last_level=1000)

# Subset multiple variables in a single dataset
ds = xr.open_mfdataset([path_to_tasmax_file, path_to_tasmin_file])
dsSub = subset_time(ds, first_level=1000.0, last_level=850.0)
```

Notes

TBA

`clisops.core.subset.subset_level_by_values` (*da*: `Union[DataArray, Dataset]`, *level_values*: `Optional[Union[Sequence[float], Sequence[int]]] = None`)
→ `Union[DataArray, Dataset]`

Subset input DataArray or Dataset based on a sequence of vertical level values.

Return a subset of a DataArray or Dataset for levels matching those requested.

Parameters

- **da** (`Union[xarray.DataArray, xarray.Dataset]`) – Input data.
- **level_values** (`Optional[Union[Sequence[float], Sequence[int]]]`) – A list of level values to select.

Returns

`Union[xarray.DataArray, xarray.Dataset]` – Subsetted `xarray.DataArray` or `xarray.Dataset`

Examples

```
import xarray as xr
from clisops.core.subset import subset_level_by_values

ds = xr.open_dataset(path_to_pr_file)

# Subset a selection of levels
levels = [1000., 850., 250., 100.]
prSub = subset_level_by_values(ds.pr, level_values=levels)
```


Notes

If any levels are not found, a `ValueError` will be raised. The requested levels will automatically be re-ordered to match the order in the input dataset.

```
clisops.core.subset.subset_shape(ds: Union[DataArray, Dataset], shape: Union[str, Path, GeoDataFrame],
                                raster_crs: Optional[Union[int, str]] = None, shape_crs:
                                Optional[Union[int, str]] = None, buffer: Optional[Union[int, float]] =
                                None, start_date: Optional[str] = None, end_date: Optional[str] = None,
                                first_level: Optional[Union[int, float]] = None, last_level:
                                Optional[Union[int, float]] = None) → Union[DataArray, Dataset]
```

Subset a `DataArray` or `Dataset` spatially (and temporally) using a vector shape and date selection.

Return a subset of a `DataArray` or `Dataset` for grid points falling within the area of a `Polygon` and/or `MultiPolygon` shape, or grid points along the path of a `LineString` and/or `MultiLineString`. If the shape consists of several disjoint polygons, the output is cut to the smallest `bbox` including all polygons.

Parameters

- **ds** (*Union[xarray.DataArray, xarray.Dataset]*) – Input values.
- **shape** (*Union[str, Path, gpd.GeoDataFrame]*) – Path to a shape file, or `GeoDataFrame` directly. Supports `GeoPandas`-compatible formats.
- **raster_crs** (*Optional[Union[str, int]]*) – EPSG number or PROJ4 string.
- **shape_crs** (*Optional[Union[str, int]]*) – EPSG number or PROJ4 string.
- **buffer** (*Optional[Union[int, float]]*) – Buffer the shape in order to select a larger region stemming from it. Units are based on the shape degrees/metres.
- **start_date** (*Optional[str]*) – Start date of the subset. Date string format – can be year (“%Y”), year-month (“%Y-%m”) or year-month-day (“%Y-%m-%d”). Defaults to first day of input data-array.
- **end_date** (*Optional[str]*) – End date of the subset. Date string format – can be year (“%Y”), year-month (“%Y-%m”) or year-month-day (“%Y-%m-%d”). Defaults to last day of input data-array.
- **first_level** (*Optional[Union[int, float]]*) – First level of the subset. Can be either an integer or float. Defaults to first level of input data-array.
- **last_level** (*Optional[Union[int, float]]*) – Last level of the subset. Can be either an integer or float. Defaults to last level of input data-array.

Returns

Union[xarray.DataArray, xarray.Dataset] – A subset of *ds*

Notes

If no CRS is found in the shape provided (e.g. RFC-7946 GeoJSON, <https://en.wikipedia.org/wiki/GeoJSON>), assumes a decimal degree datum (CRS84). Be advised that EPSG:4326 and OGC:CRS84 are not identical as axis order of lat and long differs between the two (for more information, see: <https://github.com/OSGeo/gdal/issues/2035>).

Examples

```
import xarray as xr
from clisops.core.subset import subset_shape
pr = xr.open_dataset(path_to_pr_file).pr

# Subset data array by shape
prSub = subset_shape(pr, shape=path_to_shape_file)

# Subset data array by shape and single year
prSub = subset_shape(pr, shape=path_to_shape_file, start_date='1990-01-01', end_
→date='1990-12-31')

# Subset multiple variables in a single dataset
ds = xr.open_mfdataset([path_to_tasmin_file, path_to_tasmax_file])
dsSub = subset_shape(ds, shape=path_to_shape_file)
```

`clisops.core.subset.subset_time`(*da*: Union[DataArray, Dataset], *start_date*: Optional[str] = None, *end_date*: Optional[str] = None) → Union[DataArray, Dataset]

Subset input DataArray or Dataset based on start and end years.

Return a subset of a DataArray or Dataset for dates falling within the provided bounds.

Parameters

- **da** (Union[xarray.DataArray, xarray.Dataset]) – Input data.
- **start_date** (Optional[str]) – Start date of the subset. Date string format – can be year (“%Y”), year-month (“%Y-%m”) or year-month-day (“%Y-%m-%d”). Defaults to first day of input data-array.
- **end_date** (Optional[str]) – End date of the subset. Date string format – can be year (“%Y”), year-month (“%Y-%m”) or year-month-day (“%Y-%m-%d”). Defaults to last day of input data-array.

Returns

Union[xarray.DataArray, xarray.Dataset] – Subsetted xarray.DataArray or xarray.Dataset

Examples

```
import xarray as xr
from clisops.core.subset import subset_time

ds = xr.open_dataset(path_to_pr_file)

# Subset complete years
prSub = subset_time(ds.pr, start_date='1990', end_date='1999')

# Subset single complete year
prSub = subset_time(ds.pr, start_date='1990', end_date='1990')

# Subset multiple variables in a single dataset
ds = xr.open_mfdataset([path_to_tasmax_file, path_to_tasmin_file])
dsSub = subset_time(ds, start_date='1990', end_date='1999')
```

(continues on next page)

(continued from previous page)

```
# Subset with year-month precision - Example subset 1990-03-01 to 1999-08-31
→inclusively
txSub = subset_time(ds.tasmax,start_date='1990-03',end_date='1999-08')

# Subset with specific start_dates and end_dates
tnSub = subset_time(ds.tasmin,start_date='1990-03-13',end_date='1990-08-17')
```

Notes

TODO add notes about different calendar types. Avoid “%Y-%m-31”. If you want complete month use only “%Y-%m”.

`clisops.core.subset.subset_time_by_components`(*da*: *Union[DataArray, Dataset]*, *, *time_components*: *Optional[Dict] = None*) → *DataArray*

Subsets by one or more time components (year, month, day etc).

Parameters

- **da** (*Union[xarray.DataArray, xarray.Dataset]*) – Input data.
- **time_components** (*Union[Dict, None] = None*)

Returns

xarray.DataArray

Examples

```
import xarray as xr
from clisops.core.subset import subset_time_by_components

# To select all Winter months (Dec, Jan, Feb) or [12, 1, 2]:
da = xr.open_dataset(path_to_file).pr
winter_dict = {"month": [12, 1, 2]}
res = subset_time_by_components(da, time_components=winter_dict)
```

`clisops.core.subset.subset_time_by_values`(*da*: *Union[DataArray, Dataset]*, *time_values*: *Optional[Sequence[str]] = None*) → *Union[DataArray, Dataset]*

Subset input *DataArray* or *Dataset* based on a sequence of datetime strings.

Return a subset of a *DataArray* or *Dataset* for datetimes matching those requested.

Parameters

- **da** (*Union[xarray.DataArray, xarray.Dataset]*) – Input data.
- **time_values** (*Optional[Sequence[str]]*) – Values for time. Default: *None*

Returns

Union[xarray.DataArray, xarray.Dataset] – Subsetted *xarray.DataArray* or *xarray.Dataset*

Examples

```
import xarray as xr
from clisops.core.subset import subset_time_by_values

ds = xr.open_dataset(path_to_pr_file)

# Subset a selection of datetimes
times = ["2015-01-01", "2018-12-05", "2021-06-06"]
prSub = subset_time_by_values(ds.pr, time_values=times)
```

Notes

If any datetimes are not found, a `ValueError` will be raised. The requested datetimes will automatically be re-ordered to match the order in the input dataset.

4.2 Core average functionality

Average module.

```
clisops.core.average.average_over_dims(ds: Union[DataArray, Dataset], dims: Optional[Sequence[str]]
    = None, ignore_undetected_dims: bool = False) →
    Union[DataArray, Dataset]
```

Average a `DataArray` or `Dataset` over the dimensions specified.

Parameters

- **ds** (*Union[xr.DataArray, xr.Dataset]*) – Input values.
- **dims** (*Sequence[{"time", "level", "latitude", "longitude"}]*) – The dimensions over which to apply the average. If `None`, none of the dimensions are averaged over. Dimensions must be one of [{"time", "level", "latitude", "longitude"}].
- **ignore_undetected_dims** (*bool*) – If the dimensions specified are not found in the dataset, an `Exception` will be raised if set to `True`. If `False`, an exception will not be raised and the other dimensions will be averaged over. Default = `False`

Returns

Union[xr.DataArray, xr.Dataset] – New `Dataset` or `DataArray` object averaged over the indicated dimensions. The indicated dimensions will have been removed.

Examples

```
from clisops.core.average import average_over_dims

pr = xr.open_dataset(path_to_pr_file).pr

# Average data array over latitude and longitude
prAvg = average_over_dims(pr, dims=['latitude', 'longitude'], ignore_undetected_
    ↪ dims=True)
```

`clisops.core.average.average_shape(ds: Dataset, shape: Union[str, Path, GeoDataFrame], variable: Optional[Union[str, Sequence[str]]] = None) → Union[DataArray, Dataset]`

Average a DataArray or Dataset spatially using vector shapes.

Return a DataArray or Dataset averaged over each Polygon given. Requires xESMF >= 0.5.0.

Parameters

- **ds** (*xarray.Dataset*) – Input values, coordinate attributes must be CF-compliant.
- **shape** (*Union[str, Path, gpd.GeoDataFrame]*) – Path to shape file, or directly a geodataframe. Supports formats compatible with geopandas. Will be converted to EPSG:4326 if needed.
- **variable** (*Union[str, Sequence[str], None]*) – Variables to average. If None, average over all data variables.

Returns

Union[xarray.DataArray, xarray.Dataset] – *ds* spatially-averaged over the polygon(s) in *shape*. Has a new *geom* dimension corresponding to the index of the input geodataframe. Non-geometry columns of the geodataframe are copied as auxiliary coordinates.

Notes

The spatial weights are computed with ESMF, which uses corners given in lat/lon format (EPSG:4326), the input dataset *ds* must provide those. In opposition to *subset.subset_shape*, the weights computed here take partial overlaps and holes into account.

As xESMF computes the weight masks only once, skipping missing values is not really feasible. Thus, all NaNs propagate when performing the average.

Examples

```
import xarray as xr # doctest: +SKIP
from clisops.core.average import average_shape

pr = xr.open_dataset(path_to_pr_file).pr

# Average data array over shape
prAvg = average_shape(pr, shape=path_to_shape_file)

# Average multiple variables in a single dataset
ds = xr.open_mfdataset([path_to_tasmin_file, path_to_tasmax_file])
dsAvg = average_shape(ds, shape=path_to_shape_file)
```

`clisops.core.average.average_time(ds: Union[DataArray, Dataset], freq: str) → Union[DataArray, Dataset]`

Average a DataArray or Dataset over the time frequency specified.

Parameters

- **ds** (*Union[xr.DataArray, xr.Dataset]*) – Input values.
- **freq** (*str*) – The frequency to average over. One of “month”, “year”.

Returns

Union[xr.DataArray, xr.Dataset] – New Dataset or DataArray object averaged over the indicated time frequency.

Examples

```
from clisops.core.average import average_time

pr = xr.open_dataset(path_to_pr_file).pr

# Average data array over each month
prAvg = average_time(pr, freq='month')
```

4.3 Subset operation

class `clisops.ops.subset.Subset`(*ds, file_namer='standard', split_method='time:auto', output_dir=None, output_type='netcdf', **params*)

Bases: `Operation`

`clisops.ops.subset.subset`(*ds: Union[Dataset, str, Path], *, time: Optional[Union[str, Tuple[str, str], TimeParameter, Series, Interval]] = None, area: Optional[Union[str, Tuple[Union[int, float, str], Union[int, float, str], Union[int, float, str], Union[int, float, str]], AreaParameter]] = None, level: Optional[Union[str, Tuple[Union[int, float, str], Union[int, float, str]], LevelParameter, Interval]] = None, time_components: Optional[Union[str, Dict, TimeComponents, TimeComponentsParameter]] = None, output_dir: Optional[Union[str, Path]] = None, output_type='netcdf', split_method='time:auto', file_namer='standard') → List[Union[Dataset, str]]*

Subset operation.

Parameters

- **ds** (*Union[xr.Dataset, str]*)
- **time** (*Optional[Union[str, Tuple[str, str], TimeParameter, Series, Interval]] = None*),
- **area** (*str or AreaParameter or Tuple[Union[int, float, str], Union[int, float, str], Union[int, float, str], Union[int, float, str]], optional*)
- **level** (*Optional[Union[str, Tuple[Union[int, float, str], Union[int, float, str]], LevelParameter, Interval] = None*),
- **time_components** (*Optional[Union[str, Dict, TimeComponentsParameter]] = None*),
- **output_dir** (*Optional[Union[str, Path]] = None*)
- **output_type** (*{“netcdf”, “nc”, “zarr”, “xarray”}*)
- **split_method** (*{“time:auto”}*)
- **file_namer** (*{“standard”, “simple”}*)

Returns

List[Union[xr.Dataset, str]] – A list of the subsetted outputs in the format selected; str corresponds to file paths if the output format selected is a file.

Examples

```

ds: xarray Dataset or "cmip5.output1.MOHC.HadGEM2-ES.rcp85.mon.atmos.Amon.r1i1p1.latest.tas"
time: ("1999-01-01T00:00:00", "2100-12-30T00:00:00") or "2085-01-01T12:00:00Z/2120-12-30T12:00:00Z"
area: (-5.,49.,10.,65) or "0.,49.,10.,65" or [0, 49.5, 10, 65]
level: (1000.,) or "1000/2000" or ("1000.50", "2000.60")
time_components: "year:2000,2004,2008|month:01,02" or {"year": (2000, 2004, 2008), "months": (1, 2)}
output_dir: "/cache/wps/procs/req0111"
output_type: "netcdf"
split_method: "time:auto"
file_namer: "standard"

```

Note: If you request a selection range (such as level, latitude or longitude) that specifies the lower and upper bounds in the opposite direction to the actual coordinate values then `clisops.ops.subset` will detect this issue and reverse your selection before returning the data subset.

4.4 Average operations

```

clisops.ops.average.average_over_dims(ds: Union[Dataset, str], dims: Optional[Union[Sequence[str],
    DimensionParameter]] = None, ignore_undetected_dims: bool =
    False, output_dir: Optional[Union[str, Path]] = None,
    output_type='netcdf', split_method='time:auto',
    file_namer='standard') → List[Union[Dataset, str]]

```

Parameters

- **ds** (*Union[xr.Dataset, str]*) – Xarray dataset.
- **dims** (*Optional[Union[Sequence[{"time", "level", "latitude", "longitude"}], DimensionParameter]]*) – The dimensions over which to apply the average. If `None`, none of the dimensions are averaged over. Dimensions must be one of [{"time", "level", "latitude", "longitude"}].
- **ignore_undetected_dims** (*bool*) – If the dimensions specified are not found in the dataset, an `Exception` will be raised if set to `True`. If `False`, an exception will not be raised and the other dimensions will be averaged over. Default = `False`
- **output_dir** (*Optional[Union[str, Path]]*)
- **output_type** (*{"netcdf", "nc", "zarr", "xarray"}*)
- **split_method** (*{"time:auto"}*)
- **file_namer** (*{"standard", "simple"}*)

Returns

- *List[Union[xr.Dataset, str]]*
- *A list of the outputs in the format selected; str corresponds to file paths if the*
- *output format selected is a file.*

Examples

```
ds: xarray Dataset or "cmip5.output1.MOHC.HadGEM2-ES.rcp85.mon.atmos.Amon.r1i1p1.latest.tas"
dims: ['latitude', 'longitude']
ignore_undetected_dims: False
output_dir: "/cache/wps/procs/req0111"
output_type: "netcdf"
split_method: "time:auto"
file_namer: "standard"
```

```
clisops.ops.average.average_time(ds, freq: str, output_dir: Optional[Union[str, Path]] = None,
                                output_type='netcdf', split_method='time:auto', file_namer='standard')
→ List[Union[Dataset, str]]
```

Parameters

- **ds** (*Union[xr.Dataset, str]*) – Xarray dataset.
- **freq** (*str*) – The frequency to average over. One of “month”, “year”.
- **output_dir** (*Optional[Union[str, Path]]*)
- **output_type** (*{“netcdf”, “nc”, “zarr”, “xarray”}*)
- **split_method** (*{“time:auto”}*)
- **file_namer** (*{“standard”, “simple”}*)

Returns

- *List[Union[xr.Dataset, str]]*
- *A list of the outputs in the format selected; str corresponds to file paths if the*
- *output format selected is a file.*

Examples

```
ds: xarray Dataset or "cmip5.output1.MOHC.HadGEM2-ES.rcp85.mon.atmos.Amon.r1i1p1.latest.tas"
dims: ['latitude', 'longitude']
ignore_undetected_dims: False
output_dir: "/cache/wps/procs/req0111"
output_type: "netcdf"
split_method: "time:auto"
file_namer: "standard"
```


4.5 Common functions

`clisops.utils.common.enable_logging()` → List[int]

`clisops.utils.common.expand_wildcards(paths: Union[str, Path])` → list
Expand the wildcards that may be present in Paths.

4.6 Output Utilities

`clisops.utils.output_utils.check_format(fmt)`

Checks requested format exists.

`clisops.utils.output_utils.filter_times_within(times, start=None, end=None)`

Takes an array of datetimes, returning a reduced array if start or end times are defined and are within the main array.

`clisops.utils.output_utils.get_chunk_length(da)`

Calculate the chunk length to use when chunking xarray datasets.

Based on memory limit provided in config and the size of the dataset.

`clisops.utils.output_utils.get_da(ds)`

Returns xr.DataArray when format of ds may be either xr.Dataset or xr.DataArray.

`clisops.utils.output_utils.get_format_extension(fmt)`

Finds the extension for the requested output format.

`clisops.utils.output_utils.get_format_writer(fmt)`

Finds the output method for the requested output format.

`clisops.utils.output_utils.get_output(ds, output_type, output_dir, namer)`

Return output after applying chunking and determining the output format and chunking.

`clisops.utils.output_utils.get_time_slices(ds: Union[Dataset, DataArray], split_method, start=None, end=None, file_size_limit: Optional[str] = None)` → List[Tuple[str, str]]

Take an xarray Dataset or DataArray, assume it can be split on the time axis into a sequence of slices. Optionally, take a start and end date to specify a sub-slice of the main time axis.

Use the prescribed file size limit to generate a list of (“YYYY-MM-DD”, “YYYY-MM-DD”) slices so that the output files do not (significantly) exceed the file size limit.

Parameters

- **ds** (*Union[xr.Dataset, xr.DataArray]*)
- **split_method**
- **start**
- **end**
- **file_size_limit** (*str*) – a string specifying “<number><units>”.

Returns

List[Tuple[str, str]]

4.7 File Namers

`class clisops.utils.file_namers.SimpleFileNamer(replace=None, extra="")`

Bases: `_BaseFileNamer`

Simple file namer class.

Generates numbered file names.

`class clisops.utils.file_namers.StandardFileNamer(replace=None, extra="")`

Bases: `SimpleFileNamer`

Standard file namer class.

Generates file names based on input dataset.

`get_file_name(ds, fmt='nc')`

Constructs file name.

`clisops.utils.file_namers.get_file_namer(name)`

Returns the correct filenamer from the provided name.

4.8 Dataset Utilities

`clisops.utils.dataset_utils.adjust_date_to_calendar(da, date, direction='backwards')`

Check that the date specified exists in the calendar type of the dataset.

If not present, changes the date a day at a time (up to a maximum of 5 times) to find a date that does exist. The direction to change the date by is indicated by 'direction'.

Parameters

- **da** (*xarray.Dataset* or *xarray.DataArray*) – The data to examine.
- **date** (*str*) – The date to check.
- **direction** (*str*) – The direction to move in days to find a date that does exist. 'backwards' means the search will go backwards in time until an existing date is found. 'forwards' means the search will go forwards in time. The default is 'backwards'.

Returns

str – The next possible existing date in the calendar of the dataset.

`clisops.utils.dataset_utils.calculate_offset(lon, first_element_value)`

Calculate the number of elements to roll the dataset by in order to have longitude from within requested bounds.

Parameters

- **lon** – Longitude coordinate of xarray dataset.
- **first_element_value** – The value of the first element of the longitude array to roll to.

`clisops.utils.dataset_utils.cf_convert_between_lon_frames(ds_in, lon_interval)`

Convert *ds* or *lon_interval* (whichever deems appropriate) to the other longitude frame, if the longitude frames do not match.

If *ds* and *lon_interval* are defined on different longitude frames ([-180, 180] and [0, 360]), this function will convert one of the input parameters to the other longitude frame, preferably the *lon_interval*. Adjusts shifted longitude frames [0-x, 360-x] in the dataset to one of the two standard longitude frames, dependent on the

specified `lon_interval`. In case of curvilinear grids featuring an additional 1D x-coordinate of the projection, this projection x-coordinate will not get converted.

Parameters

- **ds_in** (*xarray.Dataset or xarray.DataArray*) – xarray data object with defined longitude dimension.
- **lon_interval** (*tuple or list*) – length-2-tuple or -list of floats or integers denoting the bounds of the longitude interval.

Returns

Tuple(ds, lon_low, lon_high) – The `xarray.Dataset` and the bounds of the longitude interval, potentially adjusted in terms of their defined longitude frame.

`clisops.utils.dataset_utils.check_lon_alignment(ds, lon_bnds)`

Check whether the longitude subset requested is within the bounds of the dataset.

If not try to roll the dataset so that the request is. Raise an exception if rolling is not possible.

`clisops.utils.dataset_utils.detect_bounds(ds, coordinate) → Optional[str]`

Use `cf_xarray` to obtain the variable name of the requested coordinates bounds.

Parameters

- **ds** (*xarray.Dataset, xarray.DataArray*) – Dataset the coordinate bounds variable name shall be obtained from.
- **coordinate** (*str*) – Name of the coordinate variable to determine the bounds from.

Returns

str or None – Returns the variable name of the requested coordinate bounds. Returns `None` if the variable has no bounds or if they cannot be identified.

`clisops.utils.dataset_utils.detect_coordinate(ds, coord_type)`

Use `cf_xarray` to obtain the variable name of the requested coordinate.

Parameters

- **ds** (*xarray.Dataset, xarray.DataArray*) – Dataset the coordinate variable name shall be obtained from.
- **coord_type** (*str*) – Coordinate type understood by `cf-xarray`, eg. ‘lat’, ‘lon’, ...

Raises

AttributeError – Raised if the requested coordinate cannot be identified.

Returns

str – Coordinate variable name.

`clisops.utils.dataset_utils.detect_gridtype(ds, lon, lat, lon_bnds=None, lat_bnds=None)`

Detect type of the grid as one of “regular_lat_lon”, “curvilinear”, “unstructured”.

Assumes the grid description / structure follows the CF conventions.

4.9 Time Utilities

`clisops.utils.time_utils.create_time_bounds(ds, freq)`

Generate time bounds for datasets that have been temporally averaged.

Averaging frequencies supported are yearly, monthly and daily.

EXAMPLES

5.1 Subsetting

The subset operation makes use of `clisops.core.subset` to process the datasets and to set the output type and the output file names.

```
[1]: from clisops.utils import get_file
# fetch files locally or from github
tas_files = get_file([
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_200512-203011.nc",
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_203012-205511.nc",
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_205512-208011.nc",
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_208012-209912.nc",
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_209912-212411.nc",
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_212412-214911.nc",
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_214912-217411.nc",
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_217412-219911.nc",
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_219912-222411.nc",
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_222412-224911.nc",
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_224912-227411.nc",
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_227412-229911.nc",
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_229912-229912.nc"
])

o3_file = get_file("cmip6/o3_Amon_GFDL-ESM4_historical_r1i1p1f1_gr1_185001-194912.nc")

# remove previously created example file
import os
if os.path.exists("./output_001.nc"):
    os.remove("./output_001.nc")
```

```
[2]: from clisops.ops.subset import subset
import xarray as xr
```

The subset process takes several parameters:

5.1.1 Subsetting Parameters

```

ds: Union[xr.Dataset, str, Path]
time: Optional[Union[str, TimeParameter]]
area: Optional[
    Union[
        str,
        Tuple[
            Union[int, float, str],
            Union[int, float, str],
            Union[int, float, str],
            Union[int, float, str],
        ],
        AreaParameter,
    ]
]
level: Optional[
    Union[
        str, LevelParameter
    ]
]
time_components: Optional[Union[str, Dict, TimeComponentsParameter]]
output_dir: Optional[Union[str, Path]]
output_type: {"netcdf", "nc", "zarr", "xarray"}
split_method: {"time:auto"}
file_namer: {"standard"}

```

The output is a list containing the outputs in the format selected.

```
[3]: ds = xr.open_mfdataset(tas_files, use_cftime=True, combine="by_coords")
```

Output to xarray

There will only be one output for this example.

```
[4]: outputs = subset(
    ds=ds,
    time="2007-01-01T00:00:00/2200-12-30T00:00:00",
    area=(0.0, 10.0, 175.0, 90.0),
    output_type="xarray",
)

print(f"There is only {len(outputs)} output.")
outputs[0]
```

There is only 1 output.

```

/home/docs/checkouts/readthedocs.org/user_builds/clisops/envs/stable/lib/python3.8/site-
↳ packages/clisops/core/subset.py:1287: UserWarning: "start_date" has been nudged to
↳ nearest valid time step in xarray object.
    da = subset_time(da, start_date=start_date, end_date=end_date)
/home/docs/checkouts/readthedocs.org/user_builds/clisops/envs/stable/lib/python3.8/site-
↳ packages/clisops/core/subset.py:1287: UserWarning: "end_date" has been nudged to

```

(continues on next page)

(continued from previous page)

```
↳nearest valid time step in xarray object.
da = subset_time(da, start_date=start_date, end_date=end_date)
```

```
[4]: <xarray.Dataset>
Dimensions:    (lat: 1, time: 2329, bnds: 2, lon: 1)
Coordinates:
  height      float64 1.5
  * lat       (lat) float64 35.0
  * lon       (lon) float64 0.0
  * time      (time) object 2007-01-16 00:00:00 ... 2200-12-16 00:00:00
Dimensions without coordinates: bnds
Data variables:
  lat_bnds    (time, lat, bnds) float64 dask.array<chunksize=(287, 1, 2), meta=np.
↳ndarray>
  lon_bnds    (time, lon, bnds) float64 dask.array<chunksize=(287, 1, 2), meta=np.
↳ndarray>
  tas        (time, lat, lon) float32 dask.array<chunksize=(287, 1, 1), meta=np.
↳ndarray>
  time_bnds   (time, bnds) object dask.array<chunksize=(287, 2), meta=np.ndarray>
Attributes: (12/29)
  institution:      Met Office Hadley Centre, Fitzroy Road, Exeter, D...
  institute_id:    MOHC
  experiment_id:   rcp85
  source:          HadGEM2-ES (2009) atmosphere: HadGAM2 (N96L38); o...
  model_id:        HadGEM2-ES
  forcing:         GHG, SA, Oz, LU, S1, V1, BC, OC, (GHG = CO2, N20,...
  ...             ...
  title:           HadGEM2-ES model output prepared for CMIP5 RCP8.5
  parent_experiment: historical
  modeling_realm:  atmos
  realization:     1
  cmor_version:   2.5.0
  NCO:            4.7.3
```

Output to netCDF with simple namer

There is only one output as the file size is under the memory limit so does not need to be split. This example uses the simple namer which numbers output files.

```
[5]: outputs = subset(
    ds=ds,
    time="2007-01-01T00:00:00/2200-12-30T00:00:00",
    area=(0.0, 10.0, 175.0, 90.0),
    output_type="nc",
    output_dir=".",
    split_method="time:auto",
    file_namer="simple"
)
```

```
/home/docs/checkouts/readthedocs.org/user_builds/clisops/envs/stable/lib/python3.8/site-
↳packages/clisops/core/subset.py:1287: UserWarning: "start_date" has been nudged to
↳nearest valid time step in xarray object.
```

(continues on next page)

(continued from previous page)

```

da = subset_time(da, start_date=start_date, end_date=end_date)
/home/docs/checkouts/readthedocs.org/user_builds/clisops/envs/stable/lib/python3.8/site-
↳ packages/clisops/core/subset.py:1287: UserWarning: "end_date" has been nudged to
↳ nearest valid time step in xarray object.
da = subset_time(da, start_date=start_date, end_date=end_date)

```

```
[6]: # To open the file
```

```
subset_ds = xr.open_mfdataset("./output_001.nc", use_cftime=True, combine="by_coords")
subset_ds
```

```
[6]: <xarray.Dataset>
Dimensions:    (lat: 1, time: 2329, bnds: 2, lon: 1)
Coordinates:
  height      float64 ...
  * lat       (lat) float64 35.0
  * lon       (lon) float64 0.0
  * time      (time) object 2007-01-16 00:00:00 ... 2200-12-16 00:00:00
Dimensions without coordinates: bnds
Data variables:
  lat_bnds    (time, lat, bnds) float64 dask.array<chunksiz=(2329, 1, 2), meta=np.
↳ ndarray>
  lon_bnds    (time, lon, bnds) float64 dask.array<chunksiz=(2329, 1, 2), meta=np.
↳ ndarray>
  tas        (time, lat, lon) float32 dask.array<chunksiz=(2329, 1, 1), meta=np.
↳ ndarray>
  time_bnds   (time, bnds) object dask.array<chunksiz=(2329, 2), meta=np.ndarray>
Attributes: (12/29)
  institution:      Met Office Hadley Centre, Fitzroy Road, Exeter, D...
  institute_id:    MOHC
  experiment_id:   rcp85
  source:          HadGEM2-ES (2009) atmosphere: HadGAM2 (N96L38); o...
  model_id:        HadGEM2-ES
  forcing:         GHG, SA, Oz, LU, Sl, Vl, BC, OC, (GHG = CO2, N2O,...
  ...             ...
  title:          HadGEM2-ES model output prepared for CMIP5 RCP8.5
  parent_experiment: historical
  modeling_realm:  atmos
  realization:     1
  cmor_version:    2.5.0
  NCO:            4.7.3

```

Output to netCDF with standard namer

There is only one output as the file size is under the memory limit so does not need to be split. This example uses the standard namer which names output files according to the input file and how it has been subsetted.

```
[7]: outputs = subset(
    ds=ds,
    time="2007-01-01T00:00:00/2200-12-30T00:00:00",
    area=(0.0, 10.0, 175.0, 90.0),

```

(continues on next page)

(continued from previous page)

```

output_type="nc",
output_dir=".",
split_method="time:auto",
file_namer="standard"
)

```

```

/home/docs/checkouts/readthedocs.org/user_builds/clisops/envs/stable/lib/python3.8/site-
↳packages/clisops/core/subset.py:1287: UserWarning: "start_date" has been nudged to_
↳nearest valid time step in xarray object.
da = subset_time(da, start_date=start_date, end_date=end_date)
/home/docs/checkouts/readthedocs.org/user_builds/clisops/envs/stable/lib/python3.8/site-
↳packages/clisops/core/subset.py:1287: UserWarning: "end_date" has been nudged to_
↳nearest valid time step in xarray object.
da = subset_time(da, start_date=start_date, end_date=end_date)

```

Subsetting by level

```
[8]: ds = xr.open_dataset(o3_file, use_cftime=True)
```

No subsetting applied

```
[9]: result = subset(ds=ds,
output_type="xarray")
```

```
result[0].coords
```

```
[9]: Coordinates:
* lat      (lat) float64 -89.5 10.5
* lon      (lon) float64 0.625 125.6 250.6
* plev     (plev) float64 1e+05 9.25e+04 8.5e+04 7e+04 ... 1e+03 500.0 100.0
* time     (time) object 1850-01-16 12:00:00 ... 1949-12-16 12:00:00
```

Subsetting over level

```
[10]: # subsetting over pressure level (plev)

result = subset(ds=ds,
level="600/100",
output_type="xarray")

print(result[0].coords)
print(f"\nplev has been subsetted and now only has {len(result[0].coords)} values.")
```

```
Coordinates:
```

```

* lat      (lat) float64 -89.5 10.5
* lon      (lon) float64 0.625 125.6 250.6
* plev     (plev) float64 500.0 100.0
* time     (time) object 1850-01-16 12:00:00 ... 1949-12-16 12:00:00

```

(continues on next page)

(continued from previous page)

plev has been subsetted and now only has 4 values.

```
/home/docs/checkouts/readthedocs.org/user_builds/clisops/envs/stable/lib/python3.8/site-
↳packages/clisops/ops/subset.py:147: UserWarning: "first_level" has been nudged to
↳nearest valid level in xarray object.
result = subset_level(result, **kwargs)
```

Use time components

```
[11]: ds = xr.open_mfdataset(tas_files, use_cftime=True, combine="by_coords")
```

```
[12]: outputs = subset(
    ds=ds,
    time_components="year: 2010, 2020, 2030|month: 12, 1, 2",
    output_type="xarray",
)
```

```
print(f"There is only {len(outputs)} output.")
outputs[0]
```

There is only 1 output.

```
[12]: <xarray.Dataset>
Dimensions:    (lat: 2, time: 9, bnds: 2, lon: 2)
Coordinates:
  height      float64 1.5
  * lat       (lat) float64 -90.0 35.0
  * lon       (lon) float64 0.0 187.5
  * time      (time) object 2010-01-16 00:00:00 ... 2030-12-16 00:00:00
Dimensions without coordinates: bnds
Data variables:
  lat_bnds    (time, lat, bnds) float64 dask.array<chunksize=(8, 2, 2), meta=np.ndarray>
  lon_bnds    (time, lon, bnds) float64 dask.array<chunksize=(8, 2, 2), meta=np.ndarray>
  tas        (time, lat, lon) float32 dask.array<chunksize=(8, 2, 2), meta=np.ndarray>
  time_bnds   (time, bnds) object dask.array<chunksize=(8, 2), meta=np.ndarray>
Attributes: (12/29)
  institution:    Met Office Hadley Centre, Fitzroy Road, Exeter, D...
  institute_id:  MOHC
  experiment_id:  rcp85
  source:        HadGEM2-ES (2009) atmosphere: HadGAM2 (N96L38); o...
  model_id:      HadGEM2-ES
  forcing:       GHG, SA, Oz, LU, Sl, Vl, BC, OC, (GHG = CO2, N2O,...
  ...           ...
  title:        HadGEM2-ES model output prepared for CMIP5 RCP8.5
  parent_experiment:  historical
  modeling_realm:  atmos
  realization:    1
  cmor_version:   2.5.0
  NCO:           4.7.3
```

Using parameter classes

```
[13]: from roocs_utils.parameter.param_utils import (
    level_interval,
    level_series,
    time_components,
    time_interval,
    time_series,
)
```

```
[14]: ds = xr.open_mfdataset(tas_files, use_cftime=True, combine="by_coords")
```

```
[15]: outputs = subset(
    ds=ds,
    time=time_interval("2007-01-01T00:00:00", "2200-12-30T00:00:00"),
    time_components=time_components(month=["dec", "jan", "feb"]),
    output_type="xarray",
)
```

```
print(f"There is only {len(outputs)} output.")
outputs[0]
```

There is only 1 output.

```
/home/docs/checkouts/readthedocs.org/user_builds/clisops/envs/stable/lib/python3.8/site-
↳ packages/clisops/ops/subset.py:121: UserWarning: "start_date" has been nudged to
↳ nearest valid time step in xarray object.
    result = subset_time(self.ds, **kwargs)
/home/docs/checkouts/readthedocs.org/user_builds/clisops/envs/stable/lib/python3.8/site-
↳ packages/clisops/ops/subset.py:121: UserWarning: "end_date" has been nudged to nearest
↳ valid time step in xarray object.
    result = subset_time(self.ds, **kwargs)
```

```
[15]: <xarray.Dataset>
Dimensions:    (lat: 2, time: 583, bnds: 2, lon: 2)
Coordinates:
  height      float64 1.5
  * lat       (lat) float64 -90.0 35.0
  * lon       (lon) float64 0.0 187.5
  * time      (time) object 2007-01-16 00:00:00 ... 2200-12-16 00:00:00
Dimensions without coordinates: bnds
Data variables:
  lat_bnds    (time, lat, bnds) float64 dask.array<chunksiz=(71, 2, 2), meta=np.
↳ ndarray>
  lon_bnds    (time, lon, bnds) float64 dask.array<chunksiz=(71, 2, 2), meta=np.
↳ ndarray>
  tas         (time, lat, lon) float32 dask.array<chunksiz=(71, 2, 2), meta=np.ndarray>
  time_bnds   (time, bnds) object dask.array<chunksiz=(71, 2), meta=np.ndarray>
Attributes: (12/29)
  institution:    Met Office Hadley Centre, Fitzroy Road, Exeter, D...
  institute_id:   MOHC
  experiment_id:  rcp85
  source:         HadGEM2-ES (2009) atmosphere: HadGAM2 (N96L38); o...
  model_id:       HadGEM2-ES
```

(continues on next page)

(continued from previous page)

```

forcing:          GHG, SA, Oz, LU, S1, V1, BC, OC, (GHG = CO2, N2O,...
...
title:           HadGEM2-ES model output prepared for CMIP5 RCP8.5
parent_experiment: historical
modeling_realm:  atmos
realization:     1
cmor_version:    2.5.0
NCO:            4.7.3

```

5.2 Subsetting utilities

clisops comes with some utilities to perform common tasks that are either not implemented in xarray, or that are implemented but do not have the generality needed for climate science work. Here we show examples of the `clisops.core.subset` submodule.

```
[ ]: import xarray as xr
      from clisops.core import subset
      xr.set_options(display_style='html')

      import matplotlib.pyplot as plt
      plt.style.use('seaborn')
      plt.rcParams['figure.figsize'] = (13, 5)
      %matplotlib inline

```

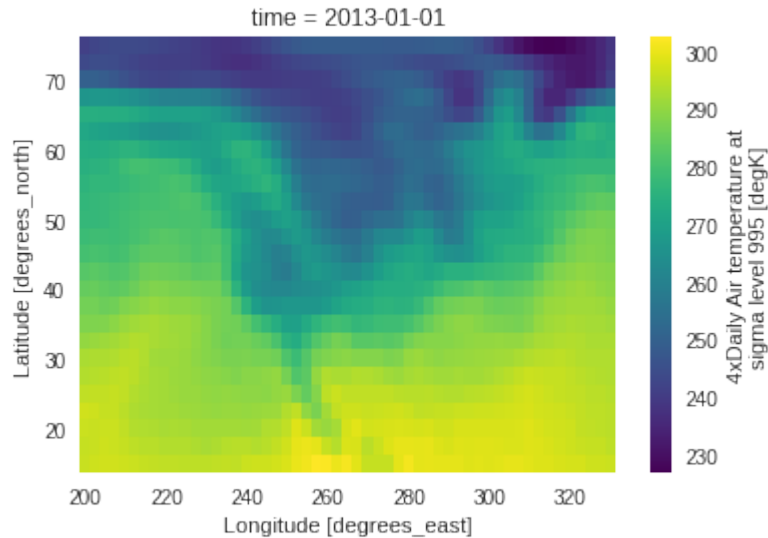
```
[ ]: ds = xr.tutorial.open_dataset('air_temperature')
      ds.coords

Coordinates:
  * lat      (lat) float32 75.0 72.5 70.0 67.5 65.0 ... 25.0 22.5 20.0 17.5 15.0
  * lon      (lon) float32 200.0 202.5 205.0 207.5 ... 322.5 325.0 327.5 330.0
  * time     (time) datetime64[ns] 2013-01-01 ... 2014-12-31T18:00:00

```

```
[ ]: ds.air.isel(time=0).plot() # Simple index-selection with xarray
<matplotlib.collections.QuadMesh at 0x7fbce02e6760>

```



nbsphinx-code-borderwhite

5.2.1 subset_bbox : using a latitude-longitude bounding box

In the previous example notebook, we used xarray's `.sel()` to cut a lat-lon subset of our data. clisops offers the same utility, but with more generality. For example, if we mindlessly try xarray's method on our dataset:

```
[ ]: ds.sel(lat=slice(45, 50), lon=slice(-60, -55)).coords
```

Coordinates:

```
* lat      (lat) float32
* lon      (lon) float32
* time     (time) datetime64[ns] 2013-01-01 ... 2014-12-31T18:00:00
```

As you can see, `lat` and `lon` are empty. In this dataset, the lats are defined in descending order and lons are in the range `[0, 360[` instead of `[-180, 180[`, which is why xarray's method did not return the expected result. clisops understands these nuances:

```
[ ]: subset.subset_bbox(ds, lat_bnds=[45, 50], lon_bnds=[-60, -55]).coords
```

Coordinates:

```
* lat      (lat) float32 50.0 47.5 45.0
* lon      (lon) float32 300.0 302.5 305.0
* time     (time) datetime64[ns] 2013-01-01 ... 2014-12-31T18:00:00
```

When lat and lon are 2-D

`subset_bbox` also manages cases where the lat-lon coordinates are not sorted 1D vectors, for example with this more complex dataset:

Most subset methods expect the input dataset / dataarray to have `lat` and `lon` as variables. It may be able to understand your data if other common names are used (like `latitude`, or `lons`), but we recommend renaming the variables before using the tool (like in this example).

```
[ ]: ds_roms = xr.tutorial.open_dataset('ROMS_example').rename(lon_rho='lon', lat_rho='lat')
salt = ds_roms.salt.isel(ocean_time=0, s_rho=0)
```

```
#matplotlib-based plotting
```

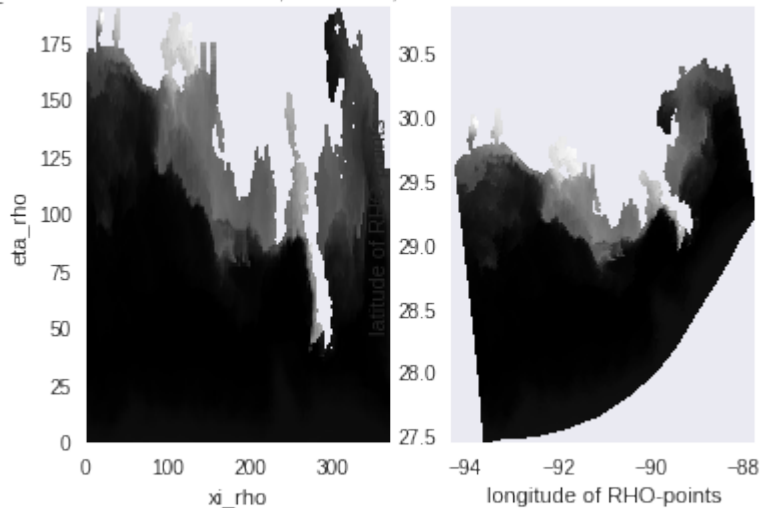
```
fig, (axEtaXi, axLatLon) = plt.subplots(1, 2)
salt.plot(cmap=plt.cm.gray_r, ax=axEtaXi, add_colorbar=False)
axLatLon.pcolormesh(salt.lon, salt.lat, salt)
axLatLon.set_xlabel(salt.lon.long_name)
axLatLon.set_ylabel(salt.lat.long_name)
```

```
<ipython-input-6-e7e8689dc693>:7: MatplotlibDeprecationWarning: shading='flat' when X,
↳and Y have the same dimensions as C is deprecated since 3.3. Either specify the
↳corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or
↳'gouraud', or set rcParams['pcolor.shading']. This will become an error two minor
↳releases later.
```

```
axLatLon.pcolormesh(salt.lon, salt.lat, salt)
```

```
Text(0, 0.5, 'latitude of RHO-points')
```

Cs_r = -0.9330103960713998, hc = 20.0, Vtransfo...



nbsphinx-code-borderwhite

```
[ ]: salt_bb = subset.subset_bbox(salt, lat_bnds=[28, 30], lon_bnds=[-91, -88])
```

```
fig, (axEtaXi, axLatLon) = plt.subplots(1, 2)
salt_bb.plot(cmap=plt.cm.gray_r, ax=axEtaXi, add_colorbar=False)
axLatLon.pcolormesh(salt_bb.lon, salt_bb.lat, salt_bb)
axLatLon.set_xlabel(salt_bb.lon.long_name)
axLatLon.set_ylabel(salt_bb.lat.long_name)
```

```
/home/phobos/Python/clisops/clisops/core/subset.py:272: UserWarning: lat and lon-like
↳dimensions are not found among arg `xarray.DataArray 'salt' (eta_rho: 191, xi_rho:
↳371)>
```

```
array([[34.953022, 34.955223, 34.957733, ..., 34.91684 , 34.913895, 34.913876],
       [34.95414 , 34.957577, 34.961086, ..., 34.916145, 34.91495 , 34.913864],
       [34.959286, 34.963936, 34.968594, ..., 34.923023, 34.92564 , 34.926853],
       ...,
       [ nan,      nan,      nan, ...,      nan,      nan,      nan],
       [ nan,      nan,      nan, ...,      nan,      nan,      nan],
```

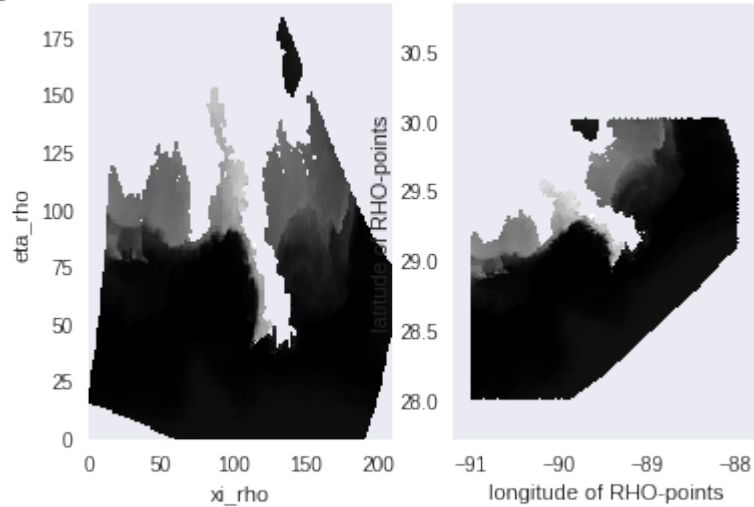
(continues on next page)

(continued from previous page)

```
[ nan, nan, nan, ..., nan, nan, nan]],
dtype=float32)
Coordinates:
Cs_r      float64 ...
lon       (eta_rho, xi_rho) float64 -93.6 -93.58 -93.57 ... -88.88 -88.87
hc        float64 ...
h         (eta_rho, xi_rho) float64 ...
lat       (eta_rho, xi_rho) float64 27.45 27.45 27.45 ... 30.85 30.86
Vtransform int32 ...
ocean_time datetime64[ns] 2001-08-01
s_rho     float64 -0.9833
Dimensions without coordinates: eta_rho, xi_rho
Attributes:
  long_name: salinity
  time:     ocean_time
  field:    salinity, scalar, series` dimensions: ['eta_rho', 'xi_rho'].
warnings.warn(
<ipython-input-7-67bc17b9a9a6>:5: MatplotlibDeprecationWarning: shading='flat' when X
↳and Y have the same dimensions as C is deprecated since 3.3. Either specify the
↳corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or
↳'gouraud', or set rcParams['pcolor.shading']. This will become an error two minor
↳releases later.
axLatLon.pcolormesh(salt_bb.lon, salt_bb.lat, salt_bb)
```

Text(0, 0.5, 'latitude of RHO-points')

Cs_r = -0.9330103960713998, hc = 20.0, Vtransfo...



nbsphinx-code-borderwhite

Add time subsetting

`subset_bbox` and other methods of the submodule give options to also give time bounds. These options are mostly for convenience as only some basics sanity checks are performed.

```
[ ]: subset.subset_bbox(ds, lat_bnds=[45, 50], lon_bnds=[-60, -55],
                       start_date='2013-02', end_date='2013-08').coords
```

Coordinates:

```
* lat      (lat) float32 50.0 47.5 45.0
* lon      (lon) float32 300.0 302.5 305.0
* time     (time) datetime64[ns] 2013-02-01 ... 2013-08-31T18:00:00
```

5.2.2 subset_gridpoint : Selecting grid points

`subset_gridpoint` can be used for selecting single locations. Compared to `.sel()`, this method adds a `tolerance` parameter in *meters* so that it finds the nearest point from the given coordinate within this distance, or else it returns NaN.

```
[ ]: lon_pt = -70.2
     lat_pt = 50.1

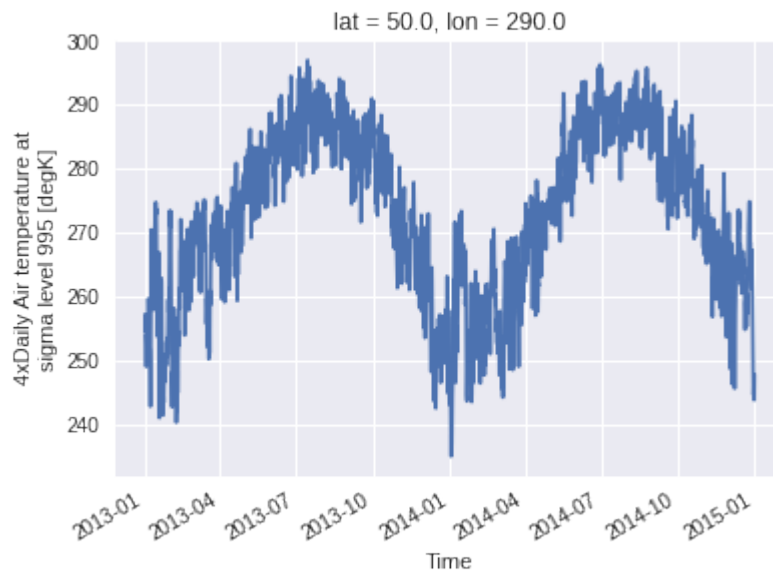
# Get timeseries for point within 100 km of (-70.2, 50.1)
ds3 = subset.subset_gridpoint(ds, lon=lon_pt, lat=lat_pt, tolerance=100 * 1000)
ds3.coords
```

Coordinates:

```
lat      float32 50.0
lon      float32 290.0
* time   (time) datetime64[ns] 2013-01-01 ... 2014-12-31T18:00:00
```

```
[ ]: ds3.air.plot()
```

```
[<matplotlib.lines.Line2D at 0x7fbcdfc62b20>]
```



nbsphinx-code-borderwhite

This also works with a list of coordinates. In that case, a new dimension `site` concatenates the selected gridpoints. Moreover, passing `add_distance = True` allows us to inspect how close the selected gridpoints were from the requested coordinates. In this example, the third point (-90, 54.1) is a little bit further than 100 km from the nearest gridpoint, so the coordinate of this gridpoint and the distance are returned, but the values are masked with NaN (and thus they do not appear on the graph).

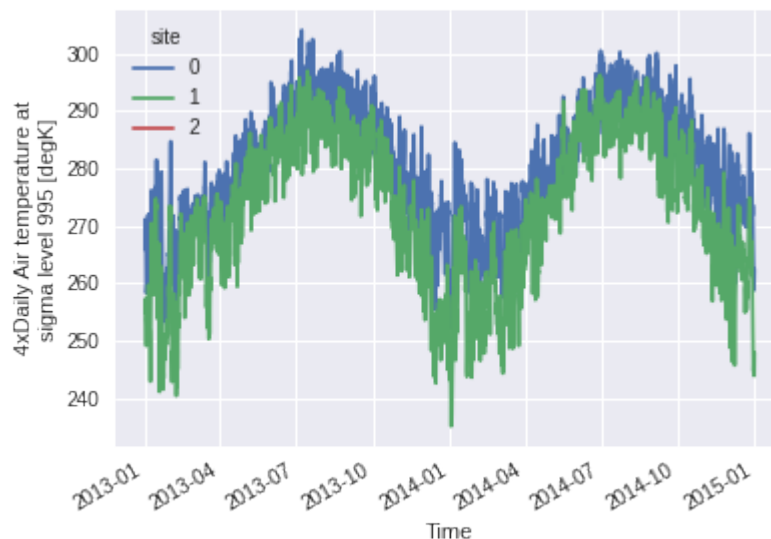
```
[ ]: lon_pt = [-65, -70.2, -90]
lat_pt = [45, 50.1, 54.1]

# Get timeseries for point within 100 km of each points
ds3 = subset.subset_gridpoint(ds, lon=lon_pt, lat=lat_pt, tolerance=100 * 1000, add_
↳distance=True)
print(ds3.coords)
ds3.air.plot(hue='site')
```

Coordinates:

```
lat      (site) float32 45.0 50.0 55.0
lon      (site) float32 295.0 290.0 270.0
* time   (time) datetime64[ns] 2013-01-01 ... 2014-12-31T18:00:00
distance (site) float64 0.0 1.814e+04 1.002e+05
```

```
[<matplotlib.lines.Line2D at 0x7fbcdfb2fa60>,
<matplotlib.lines.Line2D at 0x7fbcdfb66490>,
<matplotlib.lines.Line2D at 0x7fbcdfb66310>]
```



nbsphinx-code-borderwhite

5.2.3 subset_shape : Selecting a region from a polygon

If the region we want to keep in our dataset is complex, we can use `subset_shape` and pass a polygon. The input shape can be any georeferenced data file understood by `geopandas` and `fiona` that run under the hood. For example here with a geojson file of Canada available online:

```
[ ]: ds_can = subset.subset_shape(ds, 'https://raw.githubusercontent.com/johan/world.geo.json/
↳master/countries/CAN.geo.json')

fig, (axall, axcan) = plt.subplots(1, 2)
```

(continues on next page)

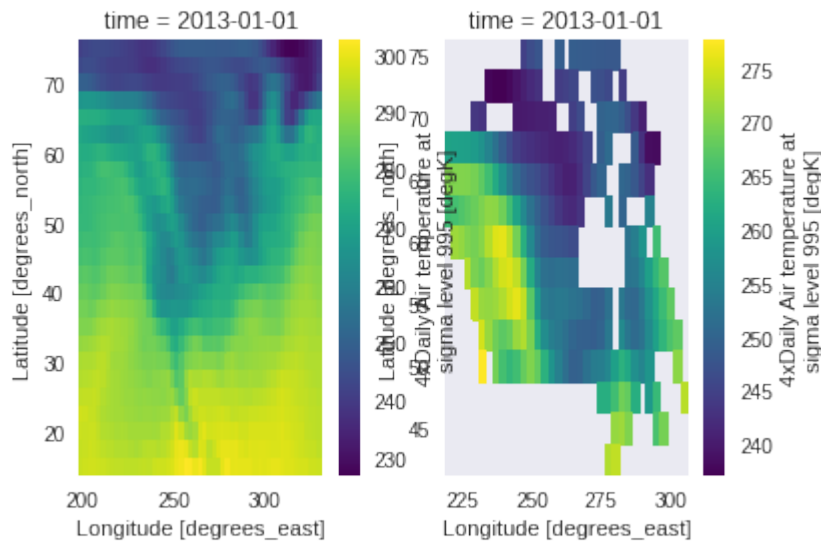
(continued from previous page)

```
ds.air.isel(time=0).plot(ax=axall)
ds_can.air.isel(time=0).plot(ax=axcan)

2021-01-26 14:44:38,120 - root - INFO - No file or no dimensions found in arg `https://
↳raw.githubusercontent.com/johan/world.geo.json/master/countries/CAN.geo.json`.

/home/phobos/Python/clisops/clisops/core/subset.py:279: UserWarning: CRS definitions are
↳similar but raster lon values must be wrapped.
    final = func(*formatted_args, **kwargs)
/home/phobos/Python/clisops/clisops/core/subset.py:463: UserWarning: Wrapping longitudes
↳at 180 degrees.
    warnings.warn("Wrapping longitudes at 180 degrees.")

<matplotlib.collections.QuadMesh at 0x7fbcdf949310>
```



nbsphinx-code-borderwhite

5.2.4 create_mask : Creating a mask from polygons

We can also create a mask to divide the data in regions. In the example, we load the shapes for Canada, Mexico and the US, merge them with pandas and then create the mask. create_mask is a bit more restrictive than the other methods showed here and specifically asks for the x and y dims. Also, we need to pass wrap_lons=True since our longitudes go from 0 to 360 instead of -180 to 180.

```
[ ]: import pandas as pd
import geopandas as gpd
can = gpd.read_file('https://raw.githubusercontent.com/johan/world.geo.json/master/
↳countries/CAN.geo.json')
usa = gpd.read_file('https://raw.githubusercontent.com/johan/world.geo.json/master/
↳countries/USA.geo.json')
mex = gpd.read_file('https://raw.githubusercontent.com/johan/world.geo.json/master/
↳countries/MEX.geo.json')
northam = pd.concat([can, usa, mex]).set_index('id')
northam
```

```
id          name \
```

(continues on next page)

(continued from previous page)

```

CAN          Canada
USA United States of America
MEX          Mexico

          geometry
id
CAN MULTIPOLYGON (((-63.66450 46.55001, -62.93930 ...
USA MULTIPOLYGON (((-155.54211 19.08348, -155.6881...
MEX POLYGON ((-97.14001 25.87000, -97.52807 24.992...

```

```
[ ]: mask = subset.create_mask(x_dim=ds.lon, y_dim=ds.lat, poly=northam, wrap_lons=True)
```

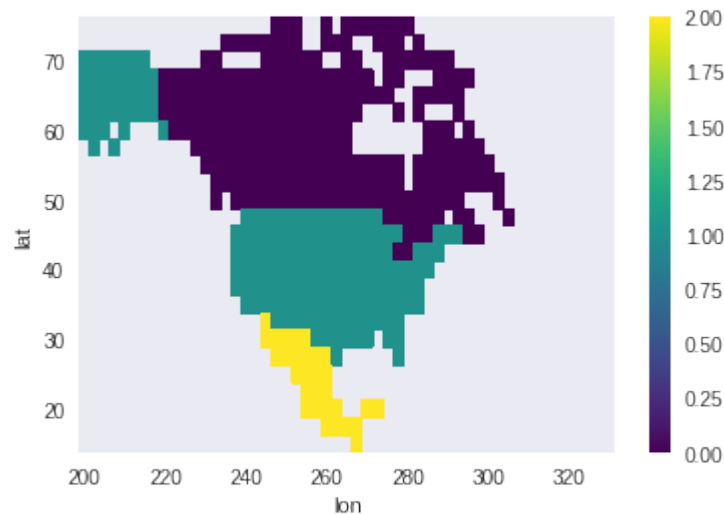
```

/home/phobos/Python/clisops/clisops/core/subset.py:463: UserWarning: Wrapping longitudes
↳ at 180 degrees.
  warnings.warn("Wrapping longitudes at 180 degrees.")

```

```
[ ]: mask.T.plot() # The transpose is needed because create_mask returns lon, lat instead the
↳ input lat, lon.
```

```
<matplotlib.collections.QuadMesh at 0x7fbcdf82ad00>
```



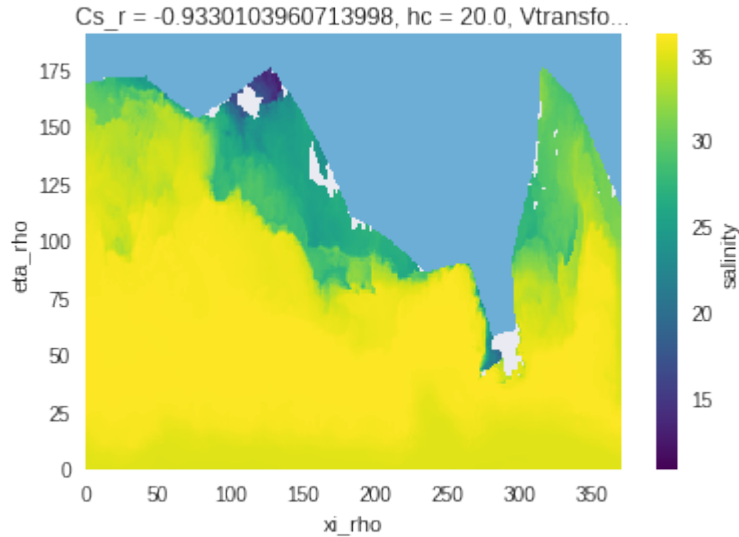
nbsphinx-code-borderwhite

The same can be done with 2D lat/lon coordinates. Here get the regions of salt that are within the boundaries of the US as defined in our geojson data. As our data comes from a small region in the Gulf of Mexico and our USA mask is quite coarse, the following isn't really helpful. In any case, it illustrates what can be done with clisops.

```
[ ]: usa_mask = subset.create_mask(x_dim=salt.lon, y_dim=salt.lat, poly=usa)
```

```
[ ]: salt.plot()
usa_mask.plot(cmap=plt.cm.Blues_r, add_colorbar=False)
```

```
<matplotlib.collections.QuadMesh at 0x7fbcdf778df0>
```



nbsphinx-code-borderwhite

5.2.5 create_weight_masks : Creating exact weights masks from polygons

While `create_mask` (and `subset_shape`) only check if the center of a grid cell lies inside of a polygon or not, `create_weight_masks` (and `average_shape`, see the Average utilities notebook) consider the partial overlap of grid cells and polygons. Thus, instead of creating an integer mask for the data, this method creates one mask for each polygon. The weights are the fraction of the grid cells-polygon overlap over the whole polygon. As such, the sum along spatial dimensions gives 1 for each geometry.

However, this feature of `clisops` requires the optional `xESMF` dependency (version 0.6.2 or later), which is not currently available of `pypi` (`pip`) and must be installed either through `conda` or from the source.

The call signature is different: a dataset or dataarray containing the grid information is passed directly. The same restrictions as for `xESMF` apply.

```
[ ]: weights = subset.create_weight_masks(ds, northam)
weights

<xarray.DataArray (geom: 3, lat: 25, lon: 53)>
array([[ [0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
          0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
        [0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
          0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
        [0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
          0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
        ...,
        [0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
          0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
        [0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
          0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
        [0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
          0.00000000e+00, 0.00000000e+00, 0.00000000e+00]],
       [[ [0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
          0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
        [0.00000000e+00, 4.38326936e-05, 1.38383820e-05, ...,
```

(continues on next page)

(continued from previous page)

```

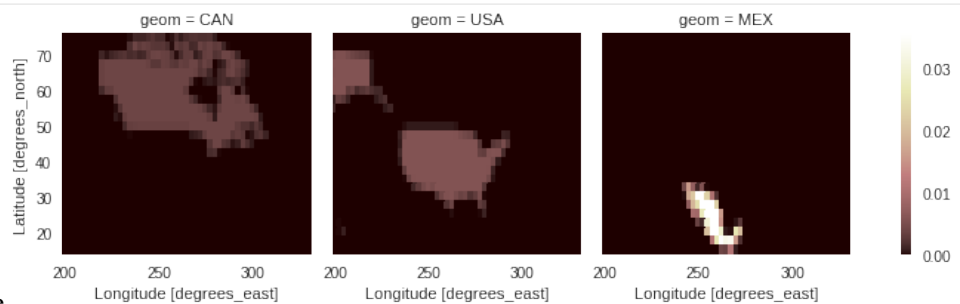
0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
[4.47662669e-03, 5.31807679e-03, 5.29956131e-03, ...,
0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
...
[0.00000000e+00, 1.73510030e-04, 9.35640399e-04, ...,
0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
[0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
[0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
0.00000000e+00, 0.00000000e+00, 0.00000000e+00]],

[[[0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
[0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
[0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
...,
[0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
[0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
0.00000000e+00, 0.00000000e+00, 0.00000000e+00],
[0.00000000e+00, 0.00000000e+00, 0.00000000e+00, ...,
0.00000000e+00, 0.00000000e+00, 0.00000000e+00]]])
Coordinates:
  name      (geom) object 'Canada' 'United States of America' 'Mexico'
* geom     (geom) object 'CAN' 'USA' 'MEX'
* lat      (lat) float32 75.0 72.5 70.0 67.5 65.0 ... 25.0 22.5 20.0 17.5 15.0
* lon      (lon) float32 200.0 202.5 205.0 207.5 ... 322.5 325.0 327.5 330.0

```

```
[ ]: weights.plot(col='geom', cmap='pink')
```

<xarray.plot.facetgrid.FacetGrid at 0x7fbcddd3f640>



nbsphinx-code-borderwhite

```
[ ]: weights.sum(['lon', 'lat'])
```

<xarray.DataArray (geom: 3)>
array([1., 1., 1.]

Coordinates:

```

  name      (geom) object 'Canada' 'United States of America' 'Mexico'
* geom     (geom) object 'CAN' 'USA' 'MEX'

```

5.3 Averaging over dimensions of the dataset

The average over dimensions operation makes use of `clisops.core.average` to process the datasets and to set the output type and the output file names.

It is possible to average over none or any number of time, longitude, latitude or level dimensions in the dataset.

```
[1]: from clisops.utils import get_file
# fetch files locally or from github
tas_files = get_file([
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_200512-203011.nc",
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_203012-205511.nc",
    "cmip5/tas_Amon_HadGEM2-ES_rcp85_r1i1p1_205512-208011.nc",
])

o3_file = get_file("cmip6/o3_Amon_GFDL-ESM4_historical_r1i1p1f1_gr1_185001-194912.nc")

# remove previously created example file
import os
if os.path.exists("./output_001.nc"):
    os.remove("./output_001.nc")
```

5.3.1 Parameters

Parameters taken by the `average_over_dims` are below:

```
ds: Union[xr.Dataset, str]
dims : Optional[Union[Tuple[str], DimensionParameter]]
    The dimensions over which to apply the average. If None, none of the dimensions are
    averaged over. Dimensions
    must be one of ["time", "level", "latitude", "longitude"].
ignore_undetected_dims: bool
    If the dimensions specified are not found in the dataset, an Exception will be raised
    if set to True.
    If False, an exception will not be raised and the other dimensions will be averaged
    over. Default = False
output_dir: Optional[Union[str, Path]] = None
output_type: {"netcdf", "nc", "zarr", "xarray"}
split_method: {"time:auto"}
file_namer: {"standard", "simple"}
```

The output is a list containing the outputs in the format selected.

```
[2]: from clisops.ops.average import average_over_dims
from roocs_utils.exceptions import InvalidParameterValue
import xarray as xr
```

```
[3]: ds = xr.open_mfdataset(tas_files, use_cftime=True, combine="by_coords")

ds
```

```
[3]: <xarray.Dataset>
Dimensions:    (lat: 2, time: 900, bnds: 2, lon: 2)
```

(continues on next page)

(continued from previous page)

```

Coordinates:
  height    float64 1.5
  * lat     (lat) float64 -90.0 35.0
  * lon     (lon) float64 0.0 187.5
  * time    (time) object 2005-12-16 00:00:00 ... 2080-11-16 00:00:00
Dimensions without coordinates: bnds
Data variables:
  lat_bnds  (time, lat, bnds) float64 dask.array<chunksize=(300, 2, 2), meta=np.
↳ndarray>
  lon_bnds  (time, lon, bnds) float64 dask.array<chunksize=(300, 2, 2), meta=np.
↳ndarray>
  tas       (time, lat, lon) float32 dask.array<chunksize=(300, 2, 2), meta=np.
↳ndarray>
  time_bnds (time, bnds) object dask.array<chunksize=(300, 2), meta=np.ndarray>
Attributes: (12/29)
  institution:      Met Office Hadley Centre, Fitzroy Road, Exeter, D...
  institute_id:    MOHC
  experiment_id:   rcp85
  source:          HadGEM2-ES (2009) atmosphere: HadGAM2 (N96L38); o...
  model_id:        HadGEM2-ES
  forcing:         GHG, SA, Oz, LU, S1, V1, BC, OC, (GHG = CO2, N2O,...
  ...              ...
  title:           HadGEM2-ES model output prepared for CMIP5 RCP8.5
  parent_experiment: historical
  modeling_realm:  atmos
  realization:     1
  cmor_version:    2.5.0
  NCO:             4.7.3

```

5.3.2 Average over one dimension

```
[4]: result = average_over_dims(ds, dims=["time"], ignore_undetected_dims=False, output_type=
↳"xarray")
```

```
result[0]
```

```
[4]: <xarray.Dataset>
Dimensions:   (lat: 2, bnds: 2, lon: 2)
Coordinates:
  height     float64 1.5
  * lat      (lat) float64 -90.0 35.0
  * lon      (lon) float64 0.0 187.5
Dimensions without coordinates: bnds
Data variables:
  lat_bnds  (lat, bnds) float64 dask.array<chunksize=(2, 2), meta=np.ndarray>
  lon_bnds  (lon, bnds) float64 dask.array<chunksize=(2, 2), meta=np.ndarray>
  tas       (lat, lon) float32 dask.array<chunksize=(2, 2), meta=np.ndarray>
Attributes: (12/29)
  institution:      Met Office Hadley Centre, Fitzroy Road, Exeter, D...
  institute_id:    MOHC

```

(continues on next page)

(continued from previous page)

```

experiment_id:      rcp85
source:             HadGEM2-ES (2009) atmosphere: HadGAM2 (N96L38); o...
model_id:          HadGEM2-ES
forcing:           GHG, SA, Oz, LU, S1, V1, BC, OC, (GHG = CO2, N20,...
...
title:             HadGEM2-ES model output prepared for CMIP5 RCP8.5
parent_experiment: historical
modeling_realm:    atmos
realization:       1
cmor_version:      2.5.0
NCO:               4.7.3

```

As you can see in the output dataset, time has been averaged over and has been removed.

5.3.3 Average over two dimensions

Averaging over two dimensions is just as simple as averaging over one. The dimensions to be averaged over should be passed in as a sequence.

```
[5]: result = average_over_dims(ds, dims=["time", "latitude"], ignore_undetected_dims=False,
→output_type="xarray")
```

```
result[0]
```

```
[5]: <xarray.Dataset>
Dimensions:  (bnds: 2, lon: 2)
Coordinates:
  height     float64 1.5
  * lon      (lon) float64 0.0 187.5
Dimensions without coordinates: bnds
Data variables:
  lat_bnds   (bnds) float64 dask.array<chunksize=(2,) , meta=np.ndarray>
  lon_bnds   (lon, bnds) float64 dask.array<chunksize=(2, 2), meta=np.ndarray>
  tas        (lon) float32 dask.array<chunksize=(2,) , meta=np.ndarray>
Attributes: (12/29)
  institution:      Met Office Hadley Centre, Fitzroy Road, Exeter, D...
  institute_id:     MOHC
  experiment_id:    rcp85
  source:           HadGEM2-ES (2009) atmosphere: HadGAM2 (N96L38); o...
  model_id:         HadGEM2-ES
  forcing:          GHG, SA, Oz, LU, S1, V1, BC, OC, (GHG = CO2, N20,...
  ...
  title:           HadGEM2-ES model output prepared for CMIP5 RCP8.5
  parent_experiment: historical
  modeling_realm:   atmos
  realization:      1
  cmor_version:     2.5.0
  NCO:              4.7.3

```

In this case both the time and latitude dimensions have been removed.

5.3.4 Allowed dimensions

It is only possible to average over longitude, latitude, level and time. If a different dimension is provided to average over an error will be raised.

```
[6]: try:
    average_over_dims(
        ds,
        dims=["incorrect_dim"],
        ignore_undetected_dims=False,
        output_type="xarray",
    )
except InvalidParameterValue as exc:
    print(exc)

Dimensions for averaging must be one of ['time', 'level', 'latitude', 'longitude']
```

5.3.5 Dimensions not found

In the case where a dimension has been selected for averaging but it doesn't exist in the dataset, there are 2 options.

1. To raise an exception when the dimension doesn't exist, set `ignore_undetected_dims = False`

```
[7]: try:
    average_over_dims(
        ds,
        dims=["level", "time"],
        ignore_undetected_dims=False,
        output_type="xarray",
    )
except InvalidParameterValue as exc:
    print(exc)

Requested dimensions were not found in input dataset: {'level'}.
```

2. To ignore when the dimension doesn't exist, and average over any other requested dimensions anyway, set `ignore_undetected_dims = True`

```
[8]: result = average_over_dims(
    ds,
    dims=["level", "time"],
    ignore_undetected_dims=True,
    output_type="xarray",
)
result[0]

[8]: <xarray.Dataset>
Dimensions:   (lat: 2, bnds: 2, lon: 2)
Coordinates:
  height      float64 1.5
  * lat       (lat) float64 -90.0 35.0
  * lon       (lon) float64 0.0 187.5
Dimensions without coordinates: bnds
Data variables:
```

(continues on next page)

(continued from previous page)

```

lat_bnds (lat, bnds) float64 dask.array<chunksize=(2, 2), meta=np.ndarray>
lon_bnds (lon, bnds) float64 dask.array<chunksize=(2, 2), meta=np.ndarray>
tas      (lat, lon) float32 dask.array<chunksize=(2, 2), meta=np.ndarray>
Attributes: (12/29)
institution:      Met Office Hadley Centre, Fitzroy Road, Exeter, D...
institute_id:    MOHC
experiment_id:   rcp85
source:          HadGEM2-ES (2009) atmosphere: HadGAM2 (N96L38); o...
model_id:        HadGEM2-ES
forcing:         GHG, SA, Oz, LU, S1, V1, BC, OC, (GHG = CO2, N2O,...
...
title:           HadGEM2-ES model output prepared for CMIP5 RCP8.5
parent_experiment: historical
modeling_realm:  atmos
realization:     1
cmor_version:    2.5.0
NCO:             4.7.3

```

In the case above, a level dimension did not exist, but this was ignored and time was averaged over anyway.

5.3.6 No dimensions supplied

If no dimensions are supplied, no averaging will be applied and the original dataset will be returned.

```

[9]: result = average_over_dims(
      ds,
      dims=None,
      ignore_undetected_dims=False,
      output_type="xarray"
    )

result[0]

-----
InvalidParameterValue                                Traceback (most recent call last)
Cell In [9], line 1
----> 1 result = average_over_dims(
      2     ds,
      3     dims=None,
      4     ignore_undetected_dims=False,
      5     output_type="xarray"
      6 )
      8 result[0]

File ~/checkouts/readthedocs.org/user_builds/clisops/envs/stable/lib/python3.8/site-
  ↳ packages/clisops/ops/average.py:89, in average_over_dims(ds, dims, ignore_undetected_
  ↳ dims, output_dir, output_type, split_method, file_namer)
      54 """
      55
      56 Parameters
      (... )
      86

```

(continues on next page)

(continued from previous page)

```

87 """
88 op = Average(**locals())
--> 89 return op.process()

```

File ~/checkouts/readthedocs.org/user_builds/clisops/envs/stable/lib/python3.8/site-
↳ packages/clisops/ops/base_operation.py:129, in Operation.process(self)

```

125 namer = self._get_file_namer()
127 # Process the xarray Dataset - this will (usually) be lazily evaluated so
128 # no actual data will be read
--> 129 processed_ds = self._calculate()
131 # remove fill values from lat/lon/time if required
132 processed_ds = self._remove_redundant_fill_values(processed_ds)

```

File ~/checkouts/readthedocs.org/user_builds/clisops/envs/stable/lib/python3.8/site-
↳ packages/clisops/ops/average.py:36, in Average._calculate(self)

```

35 def _calculate(self):
--> 36     avg_ds = average.average_over_dims(
37         self.ds,
38         self.params.get("dims", None),
39         self.params.get("ignore_undetected_dims", None),
40     )
42     return avg_ds

```

File ~/checkouts/readthedocs.org/user_builds/clisops/envs/stable/lib/python3.8/site-
↳ packages/clisops/core/average.py:182, in average_over_dims(ds, dims, ignore_undetected_
↳ dims)

```

149 """
150 Average a DataArray or Dataset over the dimensions specified.
151
152 (...)
178     prAvg = average_over_dims(pr, dims=['latitude', 'longitude'], ignore_
↳ undetected_dims=True)
179 """
181 if not dims:
--> 182     raise InvalidParameterValue(
183         "At least one dimension for averaging must be provided"
184     )
186 if not set(dims).issubset(set(known_coord_types)):
187     raise InvalidParameterValue(
188         f"Dimensions for averaging must be one of {known_coord_types}"
189     )

```

InvalidParameterValue: At least one dimension for averaging must be provided

5.3.7 An example of averaging over level

```
[10]: print("Original dataset")
print(xr.open_dataset(o3_file, use_cftime=True))

result = average_over_dims(
    o3_file,
    dims=["level"],
    ignore_undetected_dims=False,
    output_type="xarray",
)

print("Averaged dataset")
result[0]
```

```
Original dataset
<xarray.Dataset>
Dimensions:    (lat: 2, bnds: 2, lon: 3, time: 1200, plev: 19)
Coordinates:
  * lat        (lat) float64 -89.5 10.5
  * lon        (lon) float64 0.625 125.6 250.6
  * plev       (plev) float64 1e+05 9.25e+04 8.5e+04 7e+04 ... 1e+03 500.0 100.0
  * time       (time) object 1850-01-16 12:00:00 ... 1949-12-16 12:00:00
Dimensions without coordinates: bnds
Data variables:
  lat_bnds    (lat, bnds) float64 ...
  lon_bnds    (lon, bnds) float64 ...
  o3          (time, plev, lat, lon) float32 ...
  time_bnds   (time, bnds) object ...
Attributes: (12/47)
  external_variables:  areacella
  history:             Fri Oct  2 16:56:04 2020: ncks -d lat,,100 -d lo...
  table_id:           Amon
  activity_id:        CMIP
  branch_method:      standard
  branch_time_in_child: 0.0
  ...                 ...
  tracking_id:         hdl:21.14100/2601cb41-0071-4ec0-bee6-45045c81dab9
  variable_id:        o3
  variant_info:       N/A
  references:         see further_info_url attribute
  variant_label:      r1i1p1f1
  NCO:               netCDF Operators version 4.9.2 (Homepage = http://...
```

Averaged dataset

```
[10]: <xarray.Dataset>
Dimensions:    (lat: 2, lon: 3, time: 1200, bnds: 2)
Coordinates:
  * lat        (lat) float64 -89.5 10.5
  * lon        (lon) float64 0.625 125.6 250.6
  * time       (time) object 1850-01-16 12:00:00 ... 1949-12-16 12:00:00
Dimensions without coordinates: bnds
Data variables:
```

(continues on next page)

(continued from previous page)

```

o3          (time, lat, lon) float32 1.629e-06 1.63e-06 ... 1.661e-06
lat_bnds   (lat, bnds) float64 ...
lon_bnds   (lon, bnds) float64 ...
time_bnds  (time, bnds) object ...
Attributes: (12/47)
external_variables: areacella
history:      Fri Oct  2 16:56:04 2020: ncks -d lat,,,100 -d lo...
table_id:    Amon
activity_id: CMIP
branch_method: standard
branch_time_in_child: 0.0
...         ...
tracking_id: hdl:21.14100/2601cb41-0071-4ec0-bee6-45045c81dab9
variable_id: o3
variant_info: N/A
references:  see further_info_url attribute
variant_label: r1i1p1f1
NCO:        netCDF Operators version 4.9.2 (Homepage = http:/...
```

In the above, the dimension plev has be removed and averaged over

CONTRIBUTING

Contributions are welcome, and they are greatly appreciated! Every little bit helps, and credit will always be given.

You can contribute in many ways:

6.1 Types of Contributions

6.1.1 Report Bugs

Report bugs at <https://github.com/roocs/clisops/issues>.

If you are reporting a bug, please include:

- Your operating system name and version.
- Any details about your local setup that might be helpful in troubleshooting.
- Detailed steps to reproduce the bug.

6.1.2 Fix Bugs

Look through the GitHub issues for bugs. Anything tagged with “bug” and “help wanted” is open to whoever wants to implement it.

6.1.3 Implement Features

Look through the GitHub issues for features. Anything tagged with “enhancement” and “help wanted” is open to whoever wants to implement it.

6.1.4 Write Documentation

clisops could always use more documentation, whether as part of the official clisops docs, in docstrings, or even on the web in blog posts, articles, and such.

6.1.5 Submit Feedback

The best way to send feedback is to file an issue at <https://github.com/roocs/clisops/issues>.

If you are proposing a feature:

- Explain in detail how it would work.
- Keep the scope as narrow as possible, to make it easier to implement.
- Remember that this is a volunteer-driven project, and that contributions are welcome :)

6.2 Get Started!

Ready to contribute? Here's how to set up *clisops* for local development.

1. Fork the *clisops* repo on GitHub.
2. Clone your fork locally:

```
$ git clone git@github.com:roocs/clisops.git
```

3. Install your local copy into a virtualenv. Assuming you have *virtualenvwrapper* installed, this is how you set up your fork for local development:

```
# For virtualenv environments:  
$ mkvirtualenv clisops  
  
# For Anaconda/Miniconda environments:  
$ conda create -n clisops python=3.7  
  
$ cd clisops/  
$ pip install -e .
```

4. Create a branch for local development:

```
$ git checkout -b name-of-your-bugfix-or-feature
```

Now you can make your changes locally!

5. When you're done making changes, check that you verify your changes with *black* and run the tests, including testing other Python versions with *tox*:

```
# For virtualenv environments:  
$ pip install black pytest tox  
  
# For Anaconda/Miniconda environments:  
$ conda install -c conda-forge black pytest tox  
  
$ black --target-python py36 clisops tests  
$ python setup.py test  
$ tox
```

6. Before committing your changes, we ask that you install *pre-commit* in your virtualenv. *Pre-commit* runs git hooks that ensure that your code resembles that of the project and catches and corrects any small errors or inconsistencies when you *git commit*:


```
# For virtualenv environments:
$ pip install pre-commit

# For Anaconda/Miniconda environments:
$ conda install -c conda-forge pre_commit

$ pre-commit install
```

7. Commit your changes and push your branch to GitHub:

```
$ git add *

$ git commit -m "Your detailed description of your changes."
# `pre-commit` will run checks at this point:
# if no errors are found, changes will be committed.
# if errors are found, modifications will be made. Simply `git commit` again.

$ git push origin name-of-your-bugfix-or-feature
```

8. Submit a pull request through the GitHub website.

6.3 Pull Request Guidelines

Before you submit a pull request, please follow these guidelines:

1. Open an *issue* on our [GitHub repository](#) with your issue that you'd like to fix or feature that you'd like to implement.
2. Perform the changes, commit and push them either to new a branch within roocs/clisops or to your personal fork of clisops.

Warning: Try to keep your contributions within the scope of the issue that you are addressing. While it might be tempting to fix other aspects of the library as it comes up, it's better to simply to flag the problems in case others are already working on it.

Consider adding a “**# TODO:**” comment if the need arises.

3. Pull requests should raise test coverage for the clisops library. Code coverage is an indicator of how extensively tested the library is. If you are adding a new set of functions, they **must be tested** and **coverage percentage should not significantly decrease**.
4. If the pull request adds functionality, your functions should include docstring explanations. So long as the docstrings are syntactically correct, sphinx-autodoc will be able to automatically parse the information. Please ensure that the docstrings adhere to one of the following standards:
 - [numpydoc](#)
 - [reStructuredText \(ReST\)](#)
5. The pull request should work for Python 3.8+ as well as raise test coverage. Pull requests are also checked for documentation build status and for [PEP8](#) compliance.

The build statuses and build errors for pull requests can be found at:

https://travis-ci.org/roocs/clisops/pull_requests

Warning: PEP8 and Black is strongly enforced. Ensure that your changes pass **Flake8** and **Black** tests prior to pushing your final commits to your branch. Code formatting errors are treated as build errors and will block your pull request from being accepted.

6.4 Tips

To run a subset of tests:

```
$ pytest tests.test_clisops
```

6.5 Versioning

In order to update and release the library to PyPI, it's good to use a semantic versioning scheme. The method we use is as follows:

```
major.minor.patch-release
```

Major releases denote major changes resulting in a stable API;

Minor is to be used when adding a module, process or set of components;

Patch should be used for bug fixes and optimizations;

Release is a keyword used to specify the degree of production readiness (*beta* [, and optionally, *gamma*])

An increment to the Major or Minor will reset the Release to *beta*. When a build is promoted above *beta* (ie: release-ready), it's a good idea to push this version towards PyPi.

6.6 Deploying

A reminder for the maintainers on how to deploy. Make sure all your changes are committed (**including an entry in HISTORY.rst**). Then run:

```
$ bumpversion <option> # possible options: major / minor / patch / release
$ git push
$ git push --tags
```

6.7 Packaging

When test coverage and stability is adequate, maintainers should update the pip-installable package (wheel) on PyPI. In order to do this, you will need the following libraries installed:

- twine
- setuptools
- wheel

7.1 Development Lead

- Ag Stephens ag.stephens@stfc.ac.uk

7.2 Co-Developers

- Eleanor Smith eleanor.smith@stfc.ac.uk @ellesmith88
- Carsten Ehbrecht ehbrecht@dkrz.de
- Trevor James Smith smith.trevorj@ouranos.ca @Zeitsperre

7.3 Contributors

- Pascal Bourgault bourgault.pascal@ouranos.ca @aulemahal
- David Huard huard.david@ouranos.ca @huard
- Martin Schupfner schupfner@dkrz.de @sol1105

VERSION HISTORY

8.1 v0.9.3 (2022-10-03)

8.2 Bug Fixes

- Fixed a bug associated with the new `xarray` (2022.6.0+) accessor for native indexers that was introduced in (#241). (#250, #251).

8.3 Other Changes

- Fixed a handful of static type hints that were sending out warnings, despite proper use. (#251).
- Replaced all skipped doctests with sphinx-compatible python code blocks to prevent errors in downstream projects. (#251).
- Adjusted GitHub Actions builds to ensure that the `conda-xesmf` run uses the latest `xarray` available. (#251).

8.4 v0.9.2 (2022-09-06)

8.4.1 Breaking Changes

- Support has been dropped for Python3.7 and extended to Python3.10. Python3.7 is no longer tested in GitHub actions (#234).
- `packaging` has been added as a dependency (#241).

8.4.2 Bug Fixes

- Adapted `clisops.core.subset_bbox_indexer` to the newest indexing API changes in `xarray`, with backwards compatibility (#241).

8.4.3 Other Changes

- Docstrings and documentation configuration adjustments have been made to ensure that builds are adequately tested (#232, #235).

8.5 v0.9.1 (2022-05-12)

8.5.1 Bug fixes

- Fix inconsistent bounds in metadata after subset operation (#224).

8.5.2 Other Changes

- Use roocs-utils 0.6.2 to avoid test failure (#226).
- Removed unneeded testing dep from environment.yml (#223).
- Merged pre-commit autoupdate (#227).

8.6 v0.9.0 (2022-04-13)

8.6.1 New Features

- `clisops.ops.average.average_time` and `clisops.core.average.average_time` added (#211). Allowing averaging over time frequencies of day, month and year.
- New function `create_time_bounds` in `clisops.utils.time_utils`, to generate time bounds for temporally averaged datasets.
- `clisops` now uses the `loguru` library as its primary logging engine (#216). The mechanism for enabling log reporting in scripts/notebooks using `loguru` is as follows:

```
import sys
from loguru import logger

logger.activate("clisops")
LEVEL = "INFO || DEBUG || WARNING || etc."
logger.add(sys.stdout, level=LEVEL) # for logging to stdout
# or
logger.add("my_log_file.log", level=LEVEL, enqueue=True) # for logging to a file
```

8.6.2 Other Changes

- Pandas now pinned below version 1.4.0.
- Pre-commit configuration updated with code style conventions (black, pyupgrade) set to Python3.7+ (#219).
- loguru is now an install dependency, with `pytest-loguru` as a development-only dependency.
- Added function to convert the longitude axis between different longitude frames (eg. [-180, 180] and [0, 360]) (#217, #218).

8.7 v0.8.0 (2022-01-13)

8.7.1 New Features

- `clisops.core.average.average_shape` copies the global and variable attributes from the input data to the results.
- `clisops.ops.average.average_time` and `clisops.core.average.average_time` added. Allowing averaging over time frequencies of day, month and year.
- New function `create_time_bounds` in `clisops.utils.time_utils`, to generate time bounds for temporally averaged datasets.

8.7.2 Bug fixes

- `average_shape` and `create_weight_masks` were adapted to work with xESMF 0.6.2, while maintaining compatibility with earlier versions.
- Fix added to remove `_FillValue` added to coordinate variables and bounds by xarray when outputting to netCDF.

8.7.3 Other Changes

- Passing `DataArray` objects to `clisops.core.average.average_shape` is now deprecated. Averaging requires grid cell boundaries, which are not `DataArray` coordinates, but independent `Dataset` variables. Please pass `Dataset` objects and an optional list of variables to average.
- `average_shape` performs an initial subset over the averaging region, before computing the weights, to reduce memory usage.
- Minimum xesmf version set to 0.6.2.
- Minimum pygeos version set to 0.9.
- Replace `cascaded_union` by `unary_union` to anticipate a *shapely* deprecation.

8.8 v0.7.0 (2021-10-26)

8.8.1 Breaking Changes

- `time` input for `time` in `ops.subset.subset` but now be one of [`<class 'roocs_utils.parameter.param_utils.Interval'>`, `<class 'roocs_utils.parameter.param_utils.Series'>`, `<class 'NoneType'>`, `<class 'str'>`].
- `level` input for `level` in `ops.subset.subset` but now be one of [`<class 'roocs_utils.parameter.param_utils.Interval'>`, `<class 'roocs_utils.parameter.param_utils.Series'>`, `<class 'NoneType'>`, `<class 'str'>`].
- `roocs-utils` \geq 0.5.0 required.

8.8.2 New Features

- `time_values` and `level_values` arguments added to `core.subset.subset_bbox` which allows the user to provide a list of time/level values to select.
- `subset_time_by_values` and `subset_level_by_values` added to `core.subset.subset_bbox`. These allow subsetting on sequence of datetimes or levels.
- `subset_time_by_components` added to `core.subset.subset_bbox`. This allows subsetting by time components - year, month, day etc.
- `check_levels_exist` and `check_datetimes_exist` function checkers added in `core.subset` to check requested levels and datetimes exist. An exception is raised if they do not exist in the dataset.
- `time_components` argument added to `ops.subset` to allowing subsetting by time components such as year, month, day etc.

8.8.3 Other Changes

- Python 3.6 no longer tested in GitHub actions.

8.9 v0.6.5 (2021-06-10)

8.9.1 New Features

- New optional dependency `PyGEOS`, when installed the performance of `core.subset.create_mask` and `core.subset.subset_shape` are greatly improved.

8.10 v0.6.4 (2021-05-17)

8.10.1 Breaking Changes

- Exception raised in `core.average.average_over_dims` when `dims` is `None`.
- Exception raised in `core.average.average_over_shape` when `grid` and `polygon` have no overlapping values.

8.10.2 New Features

- `ops.subset.subset` now ensures all latitude and longitude bounds are in ascending order before passing to `core.subset.subset_bbox`
- `core.subset.subset_level` now checks that the order of the bounds matches the order of the level data.
- `core.subset._check_desc_coords` now checks the bounds provided are ascending before flipping them.

8.10.3 Other Changes

- `clisops` logging no longer disables other loggers.
- GitHub CI now leverages `tox` for testing as well as tests averaging functions via a conda-based build.
- Added a CI build to run against `xarray@master` that is allowed to fail.

8.11 v0.6.3 (2021-03-30)

8.11.1 Breaking Changes

- Raise an exception in `core.subset.subset_bbox` when there are no data points in the result.
- `roocs-utils` $\geq 0.3.0$ required.

8.11.2 Bug Fixes

- In `core.subset.check_start_end_dates` check if start and end date requested exist in the calendar of the dataset. If not, nudge the date forward if start date or backwards if end date.

8.11.3 Other Changes

- Error message improved to include longitude bounds of the dataset when the bounds requested in `ops.subset.subset` are not within range and rolling could not be completed.

8.12 v0.6.2 (2021-03-22)

8.12.1 Bug Fixes

- Better support for disjoint shapes in `subset_shape`.
- Identify latitude and longitude using `cf-xarray` rather than by “lat” and “lon”

8.12.2 New Features

- Add `output_staging_dir` option in `etc/roocs.ini`, to write files to initially before moving them to the requested `output_dir`.
- Notebook of examples for average over dims operation added.

8.13 v0.6.1 (2021-02-23)

8.13.1 Bug Fixes

- Add `cf-xarray` as dependency. This is a dependency of `roocs-utils` $\geq 0.2.1$ so is not a breaking change.
- Remove `python-dateutil`, `fiona` and `geojson` as dependencies, no longer needed.

8.14 v0.6.0 (2021-02-22)

8.14.1 Breaking Changes

- New dev dependency: `GitPython` $== 3.1.12$
- `roocs-utils` $\geq 0.2.1$ required.

8.14.2 New Features

- `average_over_dims` added into `average.core` and `average.ops`
- New `core.average.average_shape` + `core.subset.subset_create_weight_masks`. Depends on `xESMF` $\geq 0.5.2$, which is a new optional dependency.

8.14.3 Bug Fixes

- Fixed issue where the temporal subset was ignored if level subset selected.
- Roll dataset used in subsetting when the requested longitude bounds are not within those of the dataset.
- Fixed issue with subsetting grid lon and lat coordinates that are in descending order for `core.subset.subset_bbox`.

8.14.4 Other Changes

- Changes to allow datasets without a time dimension to be processed without issues.
- Use `DatasetMapper` from `roocs-utils` to ensure all datasets are mapped to file paths correctly.
- Using file caching to gather `mini-esgf-data` test data.
- Added a dev recipe for pip installations (*`pip install clisops[dev]`*).
- Updated pre-commit and pre-commit hooks to newest versions.
- Migrated linux-based integration builds to GitHub CI.
- Added functionality to `core.subset.create_mask` so it can accept `GeoDataFrames` with non-integer indexes.
- `clisops.utils.file_namers` adjusted to allow values to be overwritten and extras to be added to the end before the file extension.

8.15 v0.5.1 (2021-01-11)

8.15.1 Breaking Changes

- Reverting breaking changes made by the change to `core.subset.create_mask`. This change introduces a second evaluation for shapes touching grid-points.

8.15.2 Other Changes

- Using file caching to gather `xclim` test data.
- Change made to `core.subset.subset_bbox._check_desc_coords` to cope with subsetting when only one latitude or longitude exists in the input dataset

8.16 v0.5.0 (2020-12-17)

8.16.1 Breaking Changes

- Moved `core.subset.create_mask_vectorize` to `core.subset.create_mask`. The old `spatial join` option was removed.
- `core.subset.subset_shape` lost its `vectorize` kwarg, as it is now default.
- `roocs-utils>0.1.5` used

8.16.2 Other Changes

- `udunits2`>=2.2 removed as a requirement to make `clisops` completely pip installable.
- `rtee` and `libspatialindex` removed as requirements, making it easier to install through pip.
- Static types updated to include missing but permitted types.
- Better handling for paths in `ops.subset` allowing windows build to be fixed.

8.17 v0.4.0 (2020-11-10)

Adding new features, updating doc strings and documentation and inclusion of static type support.

8.17.1 Breaking Changes

- `clisops` now requires `udunits2`>=2.2.
- `roocs-utils`>=0.1.4 is now required.
- `space` parameter of `clisops.ops.subset` renamed to `area`.
- `chunk_rules` parameter of `clisops.ops.subset` renamed to `split_method`.
- `filenamer` parameter of `clisops.ops.subset` renamed to `file_namer`.

8.17.2 New Features

- `subset_level` added.
- PR template.
- Config file now exists at `clisops.etc.roocs.ini`. This can be overwritten by setting the environment variable `ROOCS_CONFIG` to the file path of a config file.
- Static typing added to subset operation function.
- `info` and `debugging` are now logged rather than printed.
- Notebook of examples for subset operation added.
- `split_method` implemented to split output files by if they exceed the memory limit provided in `clisops.etc.roocs.ini` named `file_size_limit`. Currently only the `time:auto` exists which splits evenly on time ranges.
- `file_namer` implemented in `clisops.ops.subset`. This has `simple` and `standard` options. `simple` numbers output files whereas `standard` names them according to the input dataset.
- Memory usage when completing the subsetting operation is now managed using `dask` chunking. The memory limit for memory usage for this process is set in `clisops.etc.roocs.ini` under `chunk_memory_limit`.

8.17.3 Bug Fixes

- Nudging time values to nearest available in dataset to fix a bug where subsetting failed when the exact date did not exist in the dataset.

8.17.4 Other Changes

- `cfunits` dependency removed - not needed.
- `requirements.txt` and `environment.yml` synced.
- Documentation updated to include API.
- Read the docs build now tested in CI pipeline.
- `md` files changed to `rst`.
- tests now use `mini-esgf-data` by default.

8.18 v0.3.1 (2020-08-04)

8.18.1 Other Changes

- Add missing `rtree` dependency to ensure correct spatial indexing.

8.19 v0.3.0 (2020-07-23)

8.19.1 Other Changes

- Update `testdata` and `subset` module (#34).

8.20 v0.2.1 (2020-07-08)

8.20.1 Other Changes

- Fixed docs version (#25).

8.21 v0.2.0 (2020-06-19)

8.21.1 New Features

- Integration of `xclim` subset module in `clisops.core.subset`.
- Added jupyter notebook with an example for subsetting from `xclim`.

8.21.2 Other Changes

- Fixed RTD doc build.
- Updated travis CI according to xclim requirements.
- Now employing PEP8 + Black compatible autoformatting.
- Pre-commit is now used to launch code formatting inspections for local development.

8.22 v0.1.0 (2020-04-22)

- First release.

INDICES AND TABLES

- [genindex](#)
- [modindex](#)
- [search](#)

PYTHON MODULE INDEX

C

`clisops.core.average`, 16

`clisops.core.subset`, 7

A

average_over_dims() (in module *clisops.core.average*), 16
 average_shape() (in module *clisops.core.average*), 16
 average_time() (in module *clisops.core.average*), 17

C

clisops.core.average
 module, 16
 clisops.core.subset
 module, 7
 create_mask() (in module *clisops.core.subset*), 7
 create_weight_masks() (in module *clisops.core.subset*), 8

D

distance() (in module *clisops.core.subset*), 8

M

module
 clisops.core.average, 16
 clisops.core.subset, 7

S

shape_bbox_indexer() (in module *clisops.core.subset*), 9
 subset_bbox() (in module *clisops.core.subset*), 9
 subset_gridpoint() (in module *clisops.core.subset*),
 10
 subset_level() (in module *clisops.core.subset*), 11
 subset_level_by_values() (in module *clisops.core.subset*), 12
 subset_shape() (in module *clisops.core.subset*), 13
 subset_time() (in module *clisops.core.subset*), 14
 subset_time_by_components() (in module *clisops.core.subset*), 15
 subset_time_by_values() (in module *clisops.core.subset*), 15