
nctoolkit

Robert Wilson

Mar 10, 2021

QUICK OVERVIEW

1	Installation	3
2	Introduction tutorial	5
3	News	9
4	Datasets	11
5	Importing and exporting data	15
6	Temporal statistics	17
7	Subsetting data	21
8	Manipulating variables	23
9	Interpolation	27
10	Ensemble methods	31
11	Parallel processing	35
12	Random Data Hacks	37
13	Global settings	39
14	API Reference	41
15	Package info	87
	Python Module Index	89
	Index	91

nctoolkit is a comprehensive Python package for analyzing netCDF data on Linux and macOS.

Core abilities include:

- Cropping to geographic regions
- Interactive plotting of data
- Subsetting to specific time periods
- Calculating time averages
- Calculating spatial averages
- Calculating rolling averages
- Calculating climatologies
- Creating new variables using arithmetic operations
- Calculating anomalies
- Horizontally and vertically remapping data
- Calculating the correlations between variables
- Calculating vertical averages for the likes of oceanic data
- Calculating ensemble averages
- Calculating phenological metrics

INSTALLATION

1.1 Python dependencies

- Python (3.6 or later)
- `numpy` (1.14 or later)
- `pandas` (0.24 or later)
- `xarray` (0.14 or later)
- `netCDF4` (1.53 or later)
- `hvplot`

1.2 How to install nctoolkit

The easiest way to install the package is using conda, which will install nctoolkit and all system dependencies:

```
$ conda install -c conda-forge nctoolkit
```

nctoolkit is available from the [Python Packaging Index](#). To install nctoolkit using pip:

```
$ pip install numpy  
$ pip install nctoolkit
```

If you already have numpy installed, ignore the first line. This is only included as it will make installing some dependencies smoother. nctoolkit partly relies on cartopy for plotting. This has some additional dependencies, so you may need to follow their guide [here](#) to ensure cartopy is installed fully. If you install nctoolkit using conda, you will not need to worry about that.

If you install nctoolkit from pypi, you will need to install the system dependencies listed below.

To install the development version from GitHub:

```
$ pip install git+https://github.com/r4ecology/nctoolkit.git
```

1.3 System dependencies

There are two main system dependencies: [Climate Data Operators](#), and [NCO](#). The easiest way to install them is using conda:

```
$ conda install -c conda-forge cdo
$ conda install -c conda-forge nco
```

CDO is necessary for the package to work. NCO is an optional dependency and does not have to be installed.

If you want to install CDO from source, you can use one of the bash scripts available [here](#).

INTRODUCTION TUTORIAL

nctoolkit is designed for the efficient analysis and manipulation of netCDF files. This tutorial provides an overview of how to work with individual files.

2.1 Opening netCDF data

This tutorial will illustrate the basic usage using a dataset of average global sea surface temperature from NOAA, which is available [here](#).

nctoolkit should be imported using the nc shorthand:

```
[1]: import nctoolkit as nc  
nctoolkit is using CDO version 1.9.8
```

Reading in a dataset is straightforward:

```
[2]: ff = "sst.mon.ltm.1981-2010.nc"  
sst = nc.open_data(ff)
```

We might want to know some basic information about the file. This can be done easily. Listing the available variables can be found quickly:

The current state of the dataset can be found as follows:

```
[3]: sst.variables  
[3]: ['sst', 'valid_yr_count']
```

The months available can be found using:

```
[4]: sst.months  
[4]: [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]
```

We have 12 months available. In this case it is the monthly average temperature from 1981-2010.

2.2 Modifying datasets

Each time nctoolkit executes a command that modifies a dataset, it will generate a new NetCDF file, which becomes the current file in the dataset. Before any modification this is as follows:

```
[5]: sst.current
[5]: 'sst.mon.ltm.1981-2010.nc'
```

We have seen that there are two variables in the dataset. But we only really care about `sst`. So let's select that variable:

```
[6]: sst.select(variables = "sst")
```

We can now see that there is only one variable in the `sst` dataset

```
[7]: sst.variables
[7]: ['sst']
```

We can also that a temporary file has been created with only this variable in it

```
[8]: sst.current
[8]: '/tmp/nctoolkitibehjxqnnctoolkittmpj1d6le2g.nc'
```

We have data for 12 months. But what we might really want is an average of those values. This can be quickly calculated:

```
[9]: sst.tmean()
```

Once again a new temporary file has been generated.

```
[10]: sst.current
[10]: '/tmp/nctoolkitibehjxqnnctoolkittmpmbbax0mt.nc'
```

Do not worry about the temporary folder getting clogged up. nctoolkit cleans it up automatically.

Quick visualization of netCDF data is always a good thing. So nctoolkit provides an easy autoplot feature.

```
[11]: sst.plot()
```

Data type cannot be displayed: application/javascript, application/vnd.holoviews_load.v0+json

Data type cannot be displayed: application/javascript, application/vnd.holoviews_load.v0+json

Data type cannot be displayed: application/javascript, application/vnd.holoviews_load.v0+json

Data type cannot be displayed: application/javascript, application/vnd.holoviews_load.v0+json

Data type cannot be displayed: application/javascript, application/vnd.holoviews_load.v0+json

Data type cannot be displayed: application/javascript, application/vnd.holoviews_load.v0+json

Unable to decode time axis into full numpy.datetime64 objects, continuing using `cftime.datetime` objects instead, reason: dates out of range
 Unable to decode time axis into full numpy.datetime64 objects, continuing using `cftime.datetime` objects instead, reason: dates out of range

```
[11]: :DynamicMap      []
      :Overlay
      .Image.I       :Image      [lon,lat]      (sst)
      .Coastline.I   :Feature     [Longitude,Latitude]
```

What we have seen so far is not computationally efficient. In the code below nctoolkit has generated temporary files twice:

```
[12]: sst = nc.open_data(ff)
      sst.select(variables = "sst")
      sst.tmean()
```

We can see what went on behind the scenes by accessing history:

```
[13]: sst.history
[13]: ['cdo -L -selname,sst sst.mon.ltm.1981-2010.nc /tmp/
      ↪nctoolkitibehjxqnnctoolkittmptqwymkom.nc',
      'cdo -L -timmean /tmp/nctoolkitibehjxqnnctoolkittmptqwymkom.nc /tmp/
      ↪nctoolkitibehjxqnnctoolkittmp0j30x01r.nc']
```

nctoolkit uses CDO. You do not understand how CDO works to use nctoolkit. But one nice feature of CDO is method chaining, which works like Python's. To get it working you just need to set evaluation to lazy in nctoolkit. This means nothing is evaluated until you force it to or it has to be.

```
[14]: nc.options(lazy = True)
```

Now, let's run the code again:

```
[15]: sst = nc.open_data(ff)
      sst.select(variables = "sst")
      sst.tmean()
      sst.plot()
```

Unable to decode time axis into full numpy.datetime64 objects, continuing using `cftime.datetime` objects instead, reason: dates out of range
 Unable to decode time axis into full numpy.datetime64 objects, continuing using `cftime.datetime` objects instead, reason: dates out of range

```
[15]: :DynamicMap      []
      :Overlay
      .Image.I       :Image      [lon,lat]      (sst)
      .Coastline.I   :Feature     [Longitude,Latitude]
```

When we look at history, we now see that only one temporary file was generated:

```
[16]: sst.history
[16]: ['cdo -L -timmean -selname,sst sst.mon.ltm.1981-2010.nc /tmp/
↪nctoolkitibehjxqnnctoolkittmpm0u39jql.nc']
```

In the example, above the commands were only executed when plot was called. If we want to force commands to run we use run:

```
[17]: sst = nc.open_data(ff)
sst.select_variables("sst")
sst.mean()
sst.run()
```

3.1 Release of v0.3.1

Version 0.3.1 will be released in March 2021. This is a minor release that includes new methods, under-the-hood improvements and the removal of deprecated methods.

New methods will be introduced for identifying the first time step will specific numerical thresholds are first exceeded or fallen below etc: `first_above`, `first_below`, `last_above` and `last_below`. The thresholds are either single numbers or can come from a gridded dataset for grid-cell specific thresholds.

Methods to compare a dataset with another dataset or netCDF file have been added: `gt` and `lt`, which stand for 'greater than' and 'less than'.

Users will now be able to recycle the weights calculated when interpolating data. This can enable much faster interpolation of multiple files with the same grid.

The temporal methods replaced by `tmean` etc. have now been removed from the package. So `monthly_mean` etc. can no longer be used.

3.2 Release of v0.3.0

Version 0.3.0 was released in February 2021. This will be a major release introducing major improvements to the package.

A new method `assign` is now available for generating new variables. This replaces the `mutate` and `transmute`, which were place-holder functions in the early releases of `nctoolkit` until a proper method for creating variables was put in place. `assign` operates in the same way as the `assign` method in Pandas. Users can generate new variables using lambda functions.

A major-change in this release is that evaluation is now lazy by default. The previous default of non-lazy evaluation was designed to make life slightly easier for new users of the package, but it is probably overly annoying for users to have to set evaluation to lazy each time they use the package.

This release features a subtle shift in how datasets work, so that they have consistent list-like properties. Previously, the files in a dataset given by the ``current`` attribute could be both a str or a list, depending on whether there was one or more files in the dataset. This now always gives a list. As a result datasets in `nctoolkit` have list-like properties, with ``append`` and `remove` methods available for adding and removing files. `remove` is a new method in this release. As before datasets are iterable.

This release will also allow users to run `nctoolkit` in parallel. Previous releases allowed files in multi-file datasets to be processed in parallel. However, it was not possible to create processing chains and process files in parallel. This is now possible in version thanks to under-the-hood changes in `nctoolkit`'s code base.

Users are now able to add a configuration file, which means global settings do not need to be set in every session or in every script.

DATASETS

nctoolkit works with what it calls datasets. Each dataset is made up of or more netCDF files.

4.1 Opening datasets

There are 3 ways to create a dataset: `open_data`, `open_url` or `open_thredds`.

If the data you want to analyze that is available on your computer use `open_data`. This will accept either a path to a single file or a list of files. It will also accept wildcards. So if you wanted to open all of the files in a folder called `data` as a dataset, you could do the following:

```
data = nc.open_data("data/*.nc")
```

If you want to use data that can be downloaded from a url, just use `open_url`. This will download the netCDF files to a temporary folder, and it can then be analyzed.

If you want to analyze data that is available from a thredds server or opendap, then use `open_thredds`. The file paths should end with `.nc`.

4.2 Visualization of datasets

You can visualize the contents of a dataset using the `plot` method. Below, we will plot temperature for January and the North Atlantic:

```
data = nc.open_thredds("https://psl.noaa.gov/thredds/dodsC/Datasets/COBE/data.mon.ltm.  
→1981-2010.nc")  
data.plot()
```

Please note there may be some issues due to bugs in nctoolkit's dependencies that cause problems plotting some data types. If data does not plot, raise an issue [here](#).

4.3 Modifying datasets and lazy evaluation

nctoolkit works by performing operations and then saving the results as either a temporary file or in a file specified by the user. We can illustrate this with the following code. This will select the first time step from a file available over thredds and will plot the results.

```
data = nc.open_thredds("https://psl.noaa.gov/thredds/dodsC/Datasets/COBE/data.mon.ltm.
↳1981-2010.nc")
data.select(time = 0)
data.plot()
```

You will notice, once this is done, that the file associated with the dataset is now a temporary file.

```
data.current
```

This will happen each time nctoolkit carries out an operation. This is potentially an invitation to slow-running code. You do not want to be constantly reading and writing data. Ideally, you want a processing chain which minimizes IO. nctoolkit enables this by allowing method chaining, thanks to the method chaining of its computational back-end CDO.

Let's look at this chain of code:

```
data = nc.open_thredds("https://psl.noaa.gov/thredds/dodsC/Datasets/COBE/data.mon.ltm.
↳1981-2010.nc")
data.assign(sst = lambda x: x.sst + 273.15)
data.select(months = 1)
data.crop(lon = [-80, 20], lat = [30, 70])
data.spatial_mean()
```

What is potentially wrong with this? It carries out four operations, so we absolutely do not want to create temporary file in each step. So instead of evaluating the operations line by line nctoolkit only evaluates them either when you tell it to or it has to. So in the code example above we have told nctoolkit what to do to that dataset, but have not told it to actually do any of it.

The quickest way to evaluate everything using `run`. The code above would become:

```
data = nc.open_thredds("https://psl.noaa.gov/thredds/dodsC/Datasets/COBE/data.mon.ltm.
↳1981-2010.nc")
data.assign(sst = lambda x: x.sst + 273.15)
data.select(months = 1)
data.crop(lon = [-80, 20], lat = [30, 70])
data.spatial_mean()
data.run()
```

Evaluation is, to use the technical term, lazy within nctoolkit. It only evaluates things until it needs to or is forced to.

This allows us to create efficient processing chain where we read the input file and write to the output file with no intermediate file writing. If, in the example above, we wanted to save the output file, we could do this:

```
data = nc.open_thredds("https://psl.noaa.gov/thredds/dodsC/Datasets/COBE/data.mon.ltm.
↳1981-2010.nc")
data.select(months = 1)
data.crop(lon = [-80, 20], lat = [30, 70])
data.spatial_mean()
data.to_nc("foo.nc")
```


4.4 List-like behaviour of datasets

If you want to view the files within a dataset view the `current` attribute.

This is a list that gives the file(s) within the dataset. To make processing these files easier nctoolkit features a number of methods similar to lists.

First, datasets are iterable. So, you can loop through each element of a dataset as follows:

You can find out how many files are in a dataset, using `len`:

You can add a new file to a dataset using `append`:

This method also let you add the files from another dataset.

Similarly, you can remove files from a dataset using `remove`:

In line with typical list behaviours, you can also create empty datasets as follows:

This is particularly useful if you need to create an ensemble based on multiple files that need significant processing before being added to the dataset.

4.5 Dataset attributes

We can find out key information about a dataset using its attributes.

If we want to know the variables available in a dataset called `data`, we would do:

```
data.variables
```

If we want to know the vertical levels available in the dataset, we use the following.

```
data.levels
```

If we want to know the files in a dataset, we would do this. nctoolkit works by generating temporary files, so if you have carried out any operations, this will show a list of temporary files.

```
data.current
```

If we want to find out what times are in the dataset we do this:

```
data.times
```

If we want to find out what months are in the dataset:

```
data.months
```

If we want to find out what years are in the dataset:

```
data.years
```

We can also access the history of operations carried out on the dataset. This will show the operations carried out by nctoolkit's computational back-end CDO:

```
data.history
```


IMPORTING AND EXPORTING DATA

nctoolkit can work with data available on local file systems, urls and over thredds and OPeNDAP.

5.1 Opening single files and ensembles

If you want to import a single netCDF file as a dataset, do the following:

```
import nctoolkit as nc
data = nc.open_data(infile)
```

The *open_data* function can also import multiple files. This can be done in two ways. If we have a list of files we can do the following:

```
import nctoolkit as nc
data = nc.open_data(file_list)
```

Alternatively, *open_data* is capable of handling wildcards. So if we have a folder called data, we can import all files in it as follows:

```
import nctoolkit as nc
data = nc.open_data("data/*.nc")
```

5.2 Opening files from urls/ftp

If we want to work with a file that is available at a url or ftp, we can use the *open_url* function. This will start by downloading the file to a temporary folder, so that it can be analysed.

```
import nctoolkit as nc
data = nc.open_url(example_url)
```

5.3 Opening data available over thredds servers or OPeNDAP

If you want to work with data that is available over a thredds server or OPeNDAP, you can use the *open_thredds* method. This will require that the url ends with “.nc”.

```
import nctoolkit as nc
data = nc.open_thredds(example_url)
```

5.4 Exporting datasets

nctoolkit has a number of built in methods for exporting data to netCDF, pandas dataframes and xarray datasets.

5.5 Save as a netCDF

The method *write_nc* lets users export a dataset to a netCDF file. If you want this to be a zipped netCDF file use the *zip* method before to *write_nc*. An example of usage is as follows:

```
data = nc.open_data(infile)
data.tmean()
data.zip()
data.write_nc(outfile)
```

5.6 Convert to xarray Dataset

The method *to_xarray* lets users export a dataset to an xarray dataset. An example of usage is as follows:

```
data = nc.open_data(infile)
data.tmean()
ds = data.to_xarray()
```

5.7 Convert to pandas dataframe

The method *to_dataframe* lets users export a dataset to a pandas dataframe.

```
data = nc.open_data(infile)
data.tmean()
df = data.to_dataframe()
```

TEMPORAL STATISTICS

nctoolkit has a number of built-in methods for calculating temporal statistics, all of which are prefixed with `t`: `tmean`, `tmin`, `tmax`, `trange`, `tpercentile`, `tmedian`, `tvariance`, `tstdev` and `tcumsum`.

These methods allow you to quickly calculate temporal statistics over specified time periods using the `over` argument.

By default the methods calculate the value over all time steps available. For example the following will calculate the temporal mean:

```
import nctoolkit as nc
data = nc.open_data("sst.mon.mean.nc")
data.tmean()
```

However, you may want to calculate, for example, an annual average. To do this we use `over`. This is a list which tells the function which time periods to average over. For example, the following will calculate an annual average:

```
data.tmean(["year"])
```

If you are only averaging over one time period, as above, you can simply use a character string:

```
data.tmean("year")
```

The possible options for `over` are “day”, “month”, “year”, and “season”. In this case “day” stands for day of year, not day of month.

In the example below we are calculating the maximum value in each month of each year in the dataset.

```
data.tmax(["month", "year"])
```

6.1 Calculating rolling averages

nctoolkit has a range of methods to calculate rolling averages: `rolling_mean`, `rolling_min`, `rolling_max`, `rolling_range` and `rolling_sum`. These methods let you calculate rolling statistics over a specified time window. For example, if you had daily data and you wanted to calculate a rolling weekly mean value, you could do the following:

```
data.rolling_mean(7)
```

If you wanted to calculate a rolling weekly sum, this would do:

```
data.rolling_sum(7)
```

6.2 Calculating anomalies

nctoolkit has two methods for calculating anomalies: `annual_anomaly` and `monthly_anomaly`. Both methods require you to specify a baseline period to calculate the anomaly against. They require that you specify a baseline period showing the minimum and maximum years of the climatological period to compare against.

So, if you wanted to calculate the annual anomaly compared with a baseline period of 1950-1969, you would do this:

```
data.annual_anomaly(baseline = [1950, 1969])
```

By default, the annual anomaly is calculated as the absolute difference between the annual mean in a year and the mean across the baseline period. However, in some cases this is not suitable. Instead you might want the relative change. In that case, you would do the following:

```
data.annual_anomaly(baseline = [1950, 1969], metric = "relative")
```

You can also smooth out the anomalies, so that they are calculated on a rolling basis. The following will calculate the anomaly using a rolling window of 10 years.

```
data.annual_anomaly(baseline = [1950, 1969], window = 10)
```

Monthly anomalies are calculated in the same way:

```
data.monthly_anomaly(baseline = [1950, 1969])
```

Here the anomaly is the difference between the value in each month compared with the mean in that month during the baseline period.

6.3 Calculating climatologies

This means we can easily calculate climatologies. For example the following will calculate a seasonal climatology:

```
data.tmean("season")
```

These methods allow partial matches for the arguments, which means you do not need to remember the precise argument each time. For example, the following will also calculate a seasonal climatology:

```
data.tmean("Seas")
```

Calculating a climatological monthly mean would require the following:

```
data.tmean("month")
```

and daily would be the following:

```
data.tmean("day")
```

6.4 Calculating climatologies

This means we can easily calculate climatologies. For example the following will calculate a seasonal climatology:

```
data.tmean("season")
```

6.5 Cumulative sums

We can calculate the cumulative sum as follows:

```
data.tcumsum()
```

Please note that this can only calculate over all time periods, and does not accept an `over` argument.

SUBSETTING DATA

nctoolkit has many built in methods for subsetting data. The main method is `select`. This let's you select specific variables, years, months, seasons and timesteps.

7.1 Selecting variables

If you want to select specific variables, you would do the following:

```
data.select(variables = ["var1", "var2"])
```

If you only want to select one variable, you can do this:

```
data.select(variables = "var1")
```

7.2 Selecting years

```
If you want to select specific years, you can do the following:
```

```
data.select(years = [2000, 2001])
```

Again, if you want a single year the following will work:

```
data.select(years = 2000)
```

The `select` method allows partial matches for its arguments. So if we want to select the year 2000, the following will work:

```
data.select(year = 2000)
```

In this case we can also select a range. So the following will work:

```
data.select(years = range(2000, 2010))
```

7.3 Selecting months

You can select months in the same way as years. The following examples will all do the same thing:

```
data.select(months = [1, 2, 3, 4])
data.select(months = range(1, 5))
data.select(mon = [1, 2, 3, 4])
```

7.4 Selecting seasons

You can easily select seasons. For example if you wanted to select winter, you would do the following:

```
data.select(season = "DJF")
```

7.5 Selecting timesteps

You can select specific timesteps from a dataset in a similar manner. For example if you wanted to select the first two timesteps in a dataset the following two methods will work:

```
data.select(time = [0, 1])
data.select(time = range(0, 2))
```

7.6 Geographic subsetting

If you want to select a geographic subregion of a dataset, you can use `crop`. This method will select all data within a specific longitude/latitude box. You just need to supply the minimum longitude and latitude required. In the example below, a dataset is cropped with longitudes between -80 and 90 and latitudes between 50 and 80:

```
data.crop(lon = [-80, 90], lat = [50, 80])
```

MANIPULATING VARIABLES

8.1 Creating new variables

Variable creation in nctoolkit can be done using the `assign` method, which works in a similar way to the method available in pandas.

The `assign` method works using lambda functions. Let's say we have a dataset with a variable 'var' and we simply want to add 10 to it and call the new variable 'new'. We would do the following:

```
data.assign(new = lambda x: x.var + 10)
```

If you are unfamiliar with lambda functions, note that the `x` after lambda signifies that `x` represents the dataset in whatever comes after ':', which is the actual equation to evaluate. The `x.var` term is `var` from the dataset.

By default `assign` keeps the original variables in the dataset. However, we may only want the new variable or variables. In that case you can use the `drop` argument:

```
data.assign(new = lambda x: x.var+ 10, drop = True)
```

This results in only one variable.

Note that the `assign` method uses kwargs for the lambda functions, so `drop` can be positioned anywhere. So the following will do the same thing

```
data.assign(new = lambda x: x.var+ 10, drop = True)  
data.assign(drop = True, new = lambda x: x.var+ 10)
```

At present, `assign` requires that it is written on a single line. So avoid doing something like the following:

```
data.assign(new = lambda x: x.var+ 10,  
drop = True)
```

The `assign` method will evaluate the lambda functions sent to it for each dataset grid cell for each time step. So every part of the lambda function must evaluate to a number. So the following will work:

```
k = 273.15  
data.assign(drop = True, sst_k = lambda x: x.sst + k)
```

However, if you set `k` to a string or anything other than a number it will throw an error. For example, this will throw an error:

```
k = "273.15"  
data.assign(drop = True, sst_k = lambda x: x.sst + k)
```

8.2 Applying mathematical functions to dataset variables

As part of your lambda function you can use a number of standard mathematical functions. These all have the same names as those in numpy: `abs`, `floor`, `ceil`, `sqrt`, `exp`, `log10`, `sin`, `cos`, `tan`, `arcsin`, `arccos` and `arctan`.

For example if you wanted to calculate the ceiling of a variable you could do the following:

```
data.assign(new = lambda x: ceil(x.old))
```

An example of using logs would be the following:

```
data.assign(new = lambda x: log10(x.old+1))
```

8.3 Using spatial statistics

The `assign` method carries out its calculations in each time step, and you can access spatial statistics for each time step when generating new variables. A series of functions are available that have the same names as nctoolkit methods for spatial statistics: `spatial_mean`, `spatial_max`, `spatial_min`, `spatial_sum`, `vertical_mean`, `vertical_max`, `vertical_min`, `vertical_sum`, `zonal_mean`, `zonal_max`, `zonal_min` and `zonal_sum`.

An example of the usefulness of these functions would be if you were working with global temperature data and you wanted to map regions that are warmer than average. You could do this by working out the difference between temperature in one location and the global mean:

```
data.assign(temp_comp = lambda x: x.temperature - spatial_mean(x.temperature), drop = True)
```

You can also do comparisons. In the above case, we instead might simply want to identify regions that are hotter than the global average. In that case we can simply do this:

```
data.assign(temp_comp = lambda x: x.temperature > spatial_mean(x.temperature), drop = True)
```

Let's say we wanted to map regions which are 3 degrees hotter than average. We could that as follows:

```
data.assign(temp_comp = lambda x: x.temperature > spatial_mean(x.temperature + 3), drop = True)
```

or like this:

```
data.assign(temp_comp = lambda x: x.temperature > (spatial_mean(x.temperature)+3), drop = True)
```

Logical operators work in the standard Python way. So if we had a dataset with a variable called 'var' and we wanted to find cells with values between 1 and 10, we could do this:

```
data.assign(one2ten = lambda x: x.var > 1 & x.var < 10)
```

You can process multiple variables at once using `assign`. Variables will be created in the order given, and variables created by the first lambda function can be used by the next one, and so on. The simple example below shows how this works. First we create a `var1`, which is temperature plus 1. Then `var2`, which is `var1` plus 1. Finally, we calculate the difference between `var1` and `var2`, and this should be 1 everywhere:

```
data.assign(var1 = lambda x: x.var + 1, var2 = lambda x: x.var1 + 1, diff = lambda x:
↳ x.var2 - x.var1)
```

8.4 Functions that work with nctoolkit variables

The following functions can be used on nctoolkit variables as part of lambda functions.

Function	Description	Example
abs	Absolute value	abs(x.sst)
ceiling	Ceiling of variable	ceiling(x.sst -1)
cell_area	Area of grid-cell (m2)	cell_area(x.var)
cos	Trigonometric cosine of variable	cos(x.var)
day	Day of the month of the variable	day(x.var)
exp	Exponential of variable	exp(x.sst)
floor	Floor of variable	floor(x.sst + 8.2)
hour	Hour of the day of the variable	hour(x.var)
isnan	Is variable a missing value/NA?	isnan(x.var)
latitude	Latitude of the grid cell	latitude(x.var)
level	Vertical level of variable.	level(x.var)
log	Natural log of variable	log10(x.sst + 1)
log10	Base log10 of variable	log10(x.sst + 1)
longitude	Longitude of the grid cell	longitude(x.var)
month	Month of the variable	month(x.var)
sin	Trigonometric sine of variable	sin(x.var)
spatial_max	Spatial max of variable at time-step	spatial_max(x.var)
spatial_mean	Spatial mean of variable at time-step	spatial_mean(x.var)
spatial_min	Spatial min of variable at time-step	spatial_min(x.var)
spatial_sum	Spatial sum of variable at time-step	spatial_sum(x.var)
sqrt	Square root of variable	sqrt(x.sst + 273.15)
tan	Trigonometric tangent of variable	tan(x.var)
timestep	Time step of variable. Using Python indexing.	timestep(x.var)
year	Year of the variable	year(x.var)
zonal_max	Zonal max of variable at time-step	zonal_max(x.var)
zonal_mean	Zonal mean of variable at time-step	zonal_mean(x.var)
zonal_min	Zonal min of variable at time-step	zonal_min(x.var)
zonal_sum	Zonal sum of variable at time-step	zonal_sum(x.var)

8.5 Simple mathematical operations on variables

If you want to do simple operations like adding or subtracting numbers from the variables in datasets you can use the add, subtract, divide and multiply methods. For example if you wanted to add 10 to every variable in a dataset, you would do the following:

```
data.add(10)
```

If you wanted to multiply everything by 10, you would do this:

```
data.multiply(10)
```

These methods will also let you use other datasets or netCDF files. So, you could add the values in a dataset `data2` to a dataset called `data1` as follows:

```
data1.add(data2)
```

Please note that this will require that the datasets are structured in a way that the operation makes sense. So each dimension in the datasets will either have to be identical, with the exception of when one dataset has a single value for a dimension. So for example if `data2` above has data covering only 1 timestep, but `data1` has multiple timesteps the data from that single time step will be added to all timesteps in `data1`. But if the time steps match, then the data from the first time step in `data2` will be added to the data in the first time step in `data1`, and the same will happen with the following time steps.

8.6 Simple numerical comparisons

If you want to do something as simple as working out whether the values of the variables in a dataset are greater than zero, you can use the `compare` method. This method accepts a simple comparison formula, which follows Python conventions. For example, if you wanted to figure out if the values in a dataset were greater than zero, you would do the following:

```
data.compare(">0")
```

If you wanted to know if they were equal to zero you would do this:

```
data.compare("==0")
```

INTERPOLATION

nctoolkit features built in methods for horizontal and vertical interpolation.

9.1 Interpolating to a set of coordinates

If you want to regrid a dataset to a specified set of coordinates you can `regrid` and a pandas dataframe. The first column of the dataframe should be the longitudes and the second should be latitudes. The example below regrids a sea-surface temperature dataset to a single location with longitude -30 and latitude 50.

```
import nctoolkit as nc
import pandas as pd
data = nc.open_thredds("https://psl.noaa.gov/thredds/dodsC/Datasets/COBE2/sst.mon.
↳mean.nc")
data.select(timestep = range(0, 12))
coords = pd.DataFrame({"lon":[-30], "lat":[50]})
data.regrid(coords)
```

9.1.1 Interpolating to a regular latlon grid

If you want to interpolate to a regular latlon grid, you can use `to_latlon`. `lon` and `lat` specify the minimum and maximum longitudes and latitudes, while `res`, a 2 variable list specifies the resolution. For example, if we wanted to regrid the globe to 0.5 degree north-south by 1 degree east-west resolution, we could do the following:

```
data = nc.open_thredds("https://psl.noaa.gov/thredds/dodsC/Datasets/COBE2/sst.mon.
↳mean.nc")
data.select(timestep = 0)
data.to_latlon(lon = [-79.5, 79.5], lat = [0.75, 89.75], res = [1, 0.5])
```

9.1.2 Interpolating to another dataset's grid

If we are working with two datasets and want to put them on a common grid, we can interpolate one onto the other's grid. We can illustrate this with a dataset of global sea surface temperature. Let's start by regridding the first timestep in this dataset to a regular latlon grid covering the North Atlantic.

```
data1 = nc.open_thredds("https://psl.noaa.gov/thredds/dodsC/Datasets/COBE2/sst.mon.
↳mean.nc")
data1.select(timestep = 0)
data1.to_latlon(lon = [-79.5, 79.5], lat = [-0.75, 89.75], res = [1, 0.5])
```

We can then use this new dataset as the target grid in `regrid`. So

```
data2 = nc.open_thredds("https://psl.noaa.gov/thredds/dodsC/Datasets/COBE2/sst.mon.  
↪mean.nc")  
data2.select(timestep = 0)  
data2.regrid(data1)
```

This method will also work using netCDF files. So, if you wanted you can also use a path to a netCDF file as the target grid.

9.1.3 How to reuse the weights for regridding

Please note: recycling weights only works in the dev version, and will be included in v0.3.1, which will be released publicly in March 2021 Under the hood nctoolkit regrids data by first generating a weights file. There are situations where you will want to be able to re-use these weights. For example, if you are post-processing a large number of files one after the other. To make this easier nctoolkit let's you recycle the regridding info. This let's you interpolate using either `regrid` or `to_latlon`, but keep the regridding data for future use by `regrid`.

The example below illustrates this. First, we regrid a global dataset to a regular latlon grid covering the North Atlantic, setting the recycle argument to `True`.

```
data = nc.open_thredds("https://psl.noaa.gov/thredds/dodsC/Datasets/COBE2/sst.mon.  
↪mean.nc")  
data.select(timestep = 0)  
data.to_latlon(lon = [-79.5, 79.5], lat = [-0.75, 89.75], res = [1, 0.5], recycle =   
↪True)
```

We can then use the grid from data for regridding:

```
data1 = nc.open_thredds("https://psl.noaa.gov/thredds/dodsC/Datasets/COBE2/sst.mon.  
↪mean.nc")  
data1.select(timestep = 0)  
data1.regrid(data)
```

This, of course, requires that the grids in the datasets are consistent. If you want to access the weights and grid files generated, you can do the following:

These files are deleted either when `data` is deleted or when the Python session is existed.

9.1.4 Resampling

If you want to make data more coarse spatially, just use the `resample_grid` method. This will, for example, let you select every 2nd grid cell in a north-south and east-west direction. This is illustrated in the example below, where a dataset which has spatial resolution of 1 by 1 degrees is coarsened, so that only every 10th cell is selected in a north-south and east-west. In other words it is now a 10 degrees by 10 degrees dataset.

```
data = nc.open_thredds("https://psl.noaa.gov/thredds/dodsC/Datasets/COBE2/sst.mon.  
↪mean.nc")  
data.select(timestep = 0)  
data.resample_grid(10)
```


9.1.5 Vertical interpolation

We can carry out vertical interpolation using the `vertical_interp` method. This is particularly useful for oceanic data. This is illustrated below by interpolating ocean temperatures from NOAA's World Ocean Atlas for January to a depth of 500 metres. The `vertical_interp` method requires a `levels` argument, which is sea-depth in this case.

```
data = nc.open_thredds("https://data.nodc.noaa.gov/thredds/dodsC/ncei/woa/temperature/  
↪A5B7/1.00/woa18_A5B7_t01_01.nc")  
data.select(variables="t_an")  
data.vertical_interp(levels= [500])
```


ENSEMBLE METHODS

10.1 Merging files with different variables

This notebook will outline some general methods for doing comparisons of multiple files. We will work with two different sea surface temperature data sets from NOAA and the Met Office Hadley Centre.

```
[1]: import nctoolkit as nc
import pandas as pd
import xarray as xr
import numpy as np
```

nctoolkit is using CDO version 1.9.8

Let's start by downloading the files using wget. Uncomment the code below to do this (note: you will need to extract the HadISST dataset):

```
[2]: # ! wget ftp://ftp.cdc.noaa.gov/Datasets/COBE2/sst.mon.mean.nc
# ! wget https://www.metoffice.gov.uk/hadobs/hadisst/data/HadISST_sst.nc.gz
```

The first step is to get the data. We will start by creating two separate datasets for each file.

```
[3]: sst_noaa = nc.open_data("sst.mon.mean.nc")
sst_hadley = nc.open_data("HadISST_sst.nc")
```

We can see that both variables have sea surface temperature labelled as sst. So we will need to change that.

```
[4]: sst_noaa.variables
```

```
[4]: ['sst']
```

```
[5]: sst_hadley.variables
```

```
[5]: ['sst', 'time_bnds']
```

```
[6]: sst_noaa.rename({"sst": "noaa"})
sst_hadley.rename({"sst": "hadley"})
```

The data sets also cover different time periods, and only have overlapping between 1870 and 2018. so we will need to select those years

```
[7]: sst_noaa.select(years = range(1870, 2019))
sst_hadley.select(years = range(1870, 2019))
```

We also have a problem in that there are two horizontal grids in the Hadley Centre file. We can solve this by selecting the sst variable only

```
[8]: sst_hadley.select(variables = "hadley")
```

At this point, the datasets have the same number of time steps and months covered. However, the grids are still a bit different. So we want to unify them by regridding one dataset on to the other's grid. This can be done using `regrid`, or any grid of your choosing.

```
[9]: sst_noaa.regrid(grid = sst_hadley)
```

We now have two separate datasets. Let's create a new dataset that has both of them, and then merge them. When doing this we need to make sure NAs are treated properly. In this case Hadley Centre values not being NAs as they should be, so we need to fix that. The merge method also requires a strict matching criteria for the dates in the merging files. In this case the Hadley Centre and NOAA data sets both give monthly means, but use a different day of the month. So we will set `match` to `["year", "month"]` this will ensure there are no mis-matches

```
[10]: all_sst = nc.merge(sst_noaa, sst_hadley, match = ["year", "month"])
all_sst.set_missing([-9000, -900])
```

Let's work out what the global mean SST was over the time period. Note that this will not be totally accurate as there are some missing values here and there that might bias things.

```
[11]: all_sst.spatial_mean()
all_sst.tmean("year")
all_sst.rolling_mean(10)
```

```
[12]: all_sst.plot("noaa")
```

Data type cannot be displayed: application/javascript, application/vnd.holoviews_load.v0+json

Data type cannot be displayed: application/javascript, application/vnd.holoviews_load.v0+json

Data type cannot be displayed: application/javascript, application/vnd.holoviews_load.v0+json

Data type cannot be displayed: application/javascript, application/vnd.holoviews_load.v0+json

Data type cannot be displayed: application/javascript, application/vnd.holoviews_load.v0+json

Data type cannot be displayed: application/javascript, application/vnd.holoviews_load.v0+json

```
[12]: :DynamicMap    [variable]
      :Curve       [time]    (value)
```

We can also work out the difference between the two. Here we will work out the monthly bias per cell. Then calculate the mean global difference per year, and then calculate a rolling 10 year mean.

```
[13]: all_sst = nc.open_data([sst_noaa.current, sst_hadley.current])
      all_sst.merge(match = ["year", "month"])
      all_sst.transmute({"bias": "hadley-noaa"})
      all_sst.set_missing([-9000, - 900])
      all_sst.spatial_mean()
      all_sst.tmean("year")
      all_sst.rolling_mean(10)
      all_sst.plot("bias")

[13]: :DynamicMap    [variable]
      :Curve      [time]    (value)
```

You can see that there is a notable difference at the start of the time series.

10.2 Merging files with different times

TBC

10.3 Ensemble averaging

TBC

PARALLEL PROCESSING

nctoolkit is written to enable rapid processing and analysis of netCDF files, and this includes the ability to process in parallel. Two methods of parallel processing are available. First is the ability to carry out operations on multi-file datasets in parallel. Second is the ability to define a processing chain in nctoolkit, and then use the multiprocessing package to process files in parallel using that chain.

11.1 Parallel processing of multi-file datasets

If you have a multi-file dataset, processing the files within it in parallel is easy. All you need to is the following:

```
nc.options(cores = 6)
```

This will tell nctoolkit to process the files in multi-file datasets in parallel and to use 6 cores when doing so. You can, of course, set the number of cores as high as you want. The only thing nctoolkit will do is limit it to the number of cores on your machine.

11.2 Parallel processing using multiprocessing

A common task is taking a bunch of files in a folder, doing things to them, and then saving a modified version of each file in a new folder. We want to be able to parallelize that, and we can use the multiprocessing package in the usual way.

But first, we need to change the global settings:

```
import nctoolkit as nc
nc.options(parallel = True)
```

This tells nctoolkit that we are about to do something in parallel. This is critical because of the internal workings of nctoolkit. Behind the scenes nctoolkit is constantly creating and deleting temporary files. It manages this process by creating a safe-list, i.e. a list of files in use that should not be deleted. But if you are running in parallel, you are adding to this list in parallel, and this can cause problems. Telling nctoolkit it will be run in parallel tells it to switch to using a type of list that can be safely added to in parallel.

We can use multiprocessing to do the following: take all of the files in folder foo, do a bunch of things to them, then save the results in a new folder:

We start with a function giving a processing chain. There are obviously different ways of doing this, but I like to use a function that takes the input file and output file:

```
def process_chain(infile, outfile):
    data = nc.open_data(ff)
    data.assign(tos = lambda x: x.sst + 273.15)
    data.tmean()
    data.to_nc(outfile)
```

We now want to loop through all of the files in a folder, apply the function to them and then save the results in a new folder called new:

```
ensemble = nc.create_ensemble("../data/ensemble")
import multiprocessing
pool = multiprocessing.Pool(3)
for ff in ensemble:
    pool.apply_async(process_chain, [ff, ff.replace("ensemble", "new")])
pool.close()
pool.join()
```

The number 3 in this case signifies that 3 cores are to be used.

Please note that if you are working interactively or in a Jupyter notebook, it is best to reset parallel as follows once you have stopped any parallel processing:

```
nc.options(parallel = False)
```

This is because of the effects of manually terminating commands on multiprocessing lists, which nctoolkit uses when in parallel mode.

RANDOM DATA HACKS

nctoolkit features a number of useful methods to tweak data.

12.1 Shifting time

Sometimes the times in datasets are not quite what we want, and we need some way to adjust time. An example of this is when you are missing a year of data, so want to copy data from the prior year and use it. But first you would need to shift the times in that year forward by a year. You can do this with the `shift` method. This lets you shift time forward by a specified number of hours, days, months or years. You just need to supply hours, days, months or years as an argument. So, if you wanted to shift time backward by one year, you would do the following:

```
data.shift(years = -1)
```

If you wanted to shift time forward by 12 hours, this would do it:

```
data.shift(hours = 12)
```

Note: this method allows partial matches to the arguments, so you could use hour, day, month or year just as easily.

12.2 Adding cell areas to a dataset

You can add grid cell areas to a dataset as follows:

```
data.cell_area()
```

By default, this will add the cell area (in square metres) to the dataset. If you want the dataset to only include cell areas you need to set the `join` argument to `False`:

```
data.cell_area(join = False)
```

Of course, this method will only if it is possible to calculate the areas the grid cells.

12.3 Changing the format of the netCDF files in a dataset

Sometimes you will want to change the format of the files in a dataset. You can do this using the `format` method. This let's you set the format, with the following options:

- netCDF = “nc1”
- netCDF version 2 (64-bit offset) = “nc2”/”nc”
- netCDF4 (HDF5) = “nc4”
- netCDF4-classic = “nc4c”
- netCDF version 5 (64-bit data) = “nc5”

So, if you want to set the format to netCDF4, you would do the following:

```
data.format("nc4")
```

12.4 Getting rid of dimensions with only one value

Sometimes you will have a dataset that has a dimension with only one value, and you might want to get rid of that dimension. For example, you might only have one timestep and keeping it may have no value. Getting rid of that dimension can be done using the `reduce_dims` method. It works as follows:

```
data.reduce_dims()
```

GLOBAL SETTINGS

nctoolkit lets you set global settings using options.

The most important and recommended to update is to set evaluation to lazy. This can be done as follows:

```
nc.options(lazy = True)
```

This means that commands will only be evaluated when either request them to be or they need to be.

For example, in the code below the 3 specified commands will only be calculated after it is told to `run`. This cuts down on IO, and can result in significant improvements in run time. At present lazy defaults to False, but this may change in a future release of nctoolkit.

```
nc.options(lazy = True)
data.tmean()
data.crop(lat = [0, 90])
data.spatial_mean()
data.run()
```

If you are working with ensembles, you may want to change the number of cores used for processing multiple files. For example, you can process multiple files in parallel using 6 cores as follows. By default cores = 1. Most methods can run in parallel when working with multi-file datasets.

```
nc.options(cores = 6)
```

By default nctoolkit uses the OS's temporary directories when it needs to create temporary files. In most cases this is optimal. Most of the time reading and writing to temporary folders is faster. However, in some cases this may not be a good idea because you may not have enough space in the temporary folder. In this case you can change the directory used for saving temporary files as follows:

```
nc.options(temp_dir = "/foo")
```

13.1 Setting global settings using a configuration file

You may want to set some global settings either permanently or on a project level. You can do this by setting up a configuration file. This should be a plain text file called `.nctoolkitrc` or `nctoolkitrc`. It should be placed in one of two locations: your working directory or your home directory. When nctoolkit is imported, it will look first in your working directory and then in your home directory for a file called `.nctoolkitrc` or `nctoolkitrc`. It will then use the first it finds to change the global settings from the defaults.

The structure of this file is straightforward. For example, if you wanted to set evaluation to lazy and the number of cores used for processing multi-file datasets, you would have the following in your configuration file:

lazy : True

cores : 6

The files roughly follow Python dictionary syntax, with the setting and value separate by `:`. Note that unless the setting is specified in the file, the defaults will be used. If you do not provide a configuration file, nctoolkit will use the default settings.

API REFERENCE

14.1 Session options

<code>options(**kwargs)</code>	Define session options.
---------------------------------	-------------------------

14.1.1 nctoolkit.options

`nctoolkit.options(**kwargs)`

Define session options. Set the options in the session. Available options are `thread_safe` and `lazy`. Set `thread_safe = True` if hdf5 was built to be thread safe. Set `lazy = True` if you want methods to evaluate lazy by default. Set `cores = n`, if you want nctoolkit to process the individual files in multi-file datasets in parallel. Note this only applies to multi-file datasets and will not improve performance with single files. Set `temp_dir = "/foo"` if you want to change the temporary directory used by nctoolkit to save temporary files.

Parameters ****kwargs** – Define options using key, value pairs.

Examples

If you wanted to process the files in multi-file datasets in parallel with 6 cores, do the following:

```
>>> import nctoolkit as nc
>>> nc.options(cores = 6)
```

If you want to set evaluation to always be lazy do the following:

```
>>> nc.options(lazy = True)
```

If you want nctoolkit to store temporary files in a specific directory, do this:

```
>>> nc.options(temp_dir = "/foo")
```

14.2 Opening/copying data

<code>open_data([x, checks])</code>	Read netCDF data as a DataSet object
<code>open_url([x, ftp_details, wait, file_stop])</code>	Read netCDF data from a url as a DataSet object
<code>open_thredds([x, wait, checks])</code>	Read thredds data as a DataSet object
<code>DataSet.copy(self)</code>	Make a deep copy of an DataSet object

14.2.1 nctoolkit.open_data

`nctoolkit.open_data(x=[], checks=False, **kwargs)`

Read netCDF data as a DataSet object

Parameters

- **x** (*str or list*) – A string or list of netCDF files or a single url. The function will check the files exist. If x is not a list, but an iterable it will be converted to a list. If a `*.nc` style wildcard is supplied, `open_data` will use all files available. By default an empty dataset is created, ie. using `open_data()` will create an empty dataset that can then be expanded using `append`.
- **checks** (*boolean*) – Do you want basic checks to ensure cdo can read files?
- ****kwargs** (*kwargs*) – Optional arguments for internal use by `open_thredds` and `open_url`.

Returns open_data

Return type `nctoolkit.DataSet`

Examples

If you want to open a single file as a dataset, do the following:

```
>>> import nctoolkit as nc
>>> data = nc.open_data("example.nc")
```

If you want to open a list of files as a multi-file dataset, you would do something like this:

```
>>> import nctoolkit as nc
>>> data = nc.open_data(["file1.nc", "file2.nc", "file3.nc"])
```

If you wanted to open all files in a directory “data” as a multi-file dataset, you can use a wildcard:

```
>>> import nctoolkit as nc
>>> data = nc.open_data("data/*.nc")
```

14.2.2 nctoolkit.open_url

`nctoolkit.open_url` (*x=None, ftp_details=None, wait=None, file_stop=None*)

Read netCDF data from a url as a DataSet object

Parameters

- **x** (*str*) – A string with a url. Prior to processing data will be downloaded to a temp folder.
- **ftp_details** (*dict*) – A dictionary giving the user name and password combination for ftp downloads: {"user":user, "password":pass}
- **wait** (*int*) – Time to wait, in seconds, for data to download. A minimum of 3 attempts will be made to download the data.
- **file_stop** (*int*) – Time limit, in minutes, for individual attempts at downloading data. This is useful to get around download freezes.

Returns `open_url`

Return type `nctoolkit.DataSet`

Examples

If you want to open a file available over a url do the following:

```
>>> import nctoolkit as nc
>>> data = nc.open_url("http://foo.nc")
```

This will download the file as a temporary folder for use in the dataset.

14.2.3 nctoolkit.open_thredds

`nctoolkit.open_thredds` (*x=None, wait=None, checks=False*)

Read thredds data as a DataSet object

Parameters

- **x** (*str or list*) – A string or list of thredds urls, which must end with .nc.
- **checks** (*boolean*) – Do you want to check if data is available over thredds?
- **wait** (*int*) – Time to wait for thredds server to be checked. Limitless if not supplied.

Returns `open_thredds`

Return type `nctoolkit.DataSet`

Examples

If you want to open a file available over thredds or opendap, do the following:

```
>>> import nctoolkit as nc
>>> data = nc.open_thredds("http://foo.nc")
```

14.2.4 nctoolkit.DataSet.copy

`DataSet.copy(self)`

Make a deep copy of an DataSet object

14.3 Merging or analyzing multiple datasets

<code>merge(*datasets[, match])</code>	Merge datasets
<code>cor_time([x, y])</code>	Calculate the temporal correlation coefficient between two datasets This will calculate the temporal correlation coefficient, for each time step, between two datasets.
<code>cor_space([x, y])</code>	Calculate the spatial correlation coefficient between two datasets This will calculate the spatial correlation coefficient, for each time step, between two datasets.

14.3.1 nctoolkit.merge

`nctoolkit.merge(*datasets, match=['day', 'year', 'month'])`

Merge datasets

Parameters

- **datasets** (*kwargs*) – Datasets to merge.
- **match** (*list*) – Temporal matching criteria. This is a list which must be made up of a subset of day, year, month. This checks that the datasets have compatible times. For example, if you want to ensure the datasets have the same years, then use `match = ["year"]`.

14.3.2 nctoolkit.cor_time

`nctoolkit.cor_time(x=None, y=None)`

Calculate the temporal correlation coefficient between two datasets This will calculate the temporal correlation coefficient, for each time step, between two datasets. The datasets must either have the same variables or only have one variable.

Parameters

- **x** (*dataset*) – First dataset to use
- **y** (*dataset*) – Second dataset to use

14.3.3 nctoolkit.cor_space

`nctoolkit.cor_space(x=None, y=None)`

Calculate the spatial correlation coefficient between two datasets This will calculate the spatial correlation coefficient, for each time step, between two datasets. The datasets must either have the same variables or only have one variable.

Parameters

- **x** (*dataset*) – First dataset to use
- **y** (*dataset*) – Second dataset to use

14.4 Adding file(s) to a dataset

append

14.4.1 nctoolkit.append

Functions

<code>append(self[, x])</code>	Add new file(s) to a dataset.
<code>remove(self[, x])</code>	Remove file(s) from a dataset

14.5 Accessing attributes

<code>DataSet.variables</code>	List variables contained in a dataset
<code>DataSet.years</code>	List years contained in a dataset
<code>DataSet.months</code>	List months contained in a dataset
<code>DataSet.times</code>	List times contained in a dataset
<code>DataSet.levels</code>	List levels contained in a dataset
<code>DataSet.size</code>	The size of an object This will print the number of files, total size, and smallest and largest files in an DataSet object.
<code>DataSet.current</code>	The current file or files in the DataSet object
<code>DataSet.history</code>	The history of operations on the DataSet
<code>DataSet.start</code>	The starting file or files of the DataSet object

14.5.1 nctoolkit.DataSet.variables

property `DataSet.variables`
List variables contained in a dataset

14.5.2 nctoolkit.DataSet.years

property `DataSet.years`
List years contained in a dataset

14.5.3 nctoolkit.DataSet.months

property `DataSet.months`
List months contained in a dataset

14.5.4 nctoolkit.DataSet.times

property `DataSet.times`
List times contained in a dataset

14.5.5 nctoolkit.DataSet.levels

property `DataSet.levels`
List levels contained in a dataset

14.5.6 nctoolkit.DataSet.size

property `DataSet.size`
The size of an object This will print the number of files, total size, and smallest and largest files in an DataSet object.

14.5.7 nctoolkit.DataSet.current

property `DataSet.current`
The current file or files in the DataSet object

14.5.8 nctoolkit.DataSet.history

property `DataSet.history`
The history of operations on the DataSet

14.5.9 nctoolkit.DataSet.start

property `DataSet.start`
The starting file or files of the DataSet object

14.6 Plotting

```
DataSet.plot(self[, vars])
```

14.6.1 nctoolkit.DataSet.plot

```
DataSet.plot (self, vars=None)
```

14.7 Variable modification

<i>DataSet.assign</i> (self[, drop])	Create new variables Existing columns that are re-assigned will be overwritten. :param drop: Set to True if you want existing variables to be removed once the new ones have been created. Defaults to False.
<i>DataSet.rename</i> (self, newnames)	Rename variables in a dataset
<i>DataSet.set_missing</i> (self[, value])	Set the missing value for a single number or a range
<i>DataSet.sum_all</i> (self[, drop])	Calculate the sum of all variables for each time step

14.7.1 nctoolkit.DataSet.assign

```
DataSet.assign (self, drop=False, \**kwargs)
```

Create new variables Existing columns that are re-assigned will be overwritten. :param drop: Set to True if you want existing variables to be removed once the new ones have been created.

Defaults to False.

should evaluate to a numeric. New variables are calculated for each grid cell and time step.

Parameters ****kwargs** (*dict of {str: callable}*) – New variable names are keywords. All terms in the equation given by the lamda function should evaluate to a numeric. New variables are calculated for each grid cell and time step.

Notes

Operations are carried out in the order give. So if a new variable is created in the first argument, it can then be used in following arguments.

14.7.2 nctoolkit.DataSet.rename

`DataSet.rename(self, newnames)`

Rename variables in a dataset

Parameters `newnames` (*dict*) – Dictionary with key-value pairs being original and new variable names

Examples

If you want to rename a variable x to y, do the following:

```
>>> data.rename({"x": "y"})
```

14.7.3 nctoolkit.DataSet.set_missing

`DataSet.set_missing(self, value=None)`

Set the missing value for a single number or a range

Parameters `value` (*2 variable list or int/float*) – If int/float is provided, the missing value will be set to that. If a list is provided, values between the two values (inclusive) of the list are set to missing.

14.7.4 nctoolkit.DataSet.sum_all

`DataSet.sum_all(self, drop=True)`

Calculate the sum of all variables for each time step

Parameters `drop` (*boolean*) – Do you want to keep variables?

14.8 netCDF file attribute modification

<code>DataSet.set_longnames(self[, name_dict])</code>	Set the long names of variables
<code>DataSet.set_units(self[, unit_dict])</code>	Set the units for variables

14.8.1 nctoolkit.DataSet.set_longnames

`DataSet.set_longnames(self, name_dict=None)`

Set the long names of variables

Parameters `name_dict` (*dict*) – Dictionary with key, value pairs representing the variable names and their long names

14.8.2 nctoolkit.DataSet.set_units

`DataSet.set_units(self, unit_dict=None)`

Set the units for variables

Parameters `unit_dict` (*dict*) – A dictionary where the key-value pairs are the variables and new units respectively.

14.9 Vertical/level methods

<code>DataSet.surface(self)</code>	Extract the top/surface level from a dataset This extracts the first vertical level from each file in a dataset.
<code>DataSet.bottom(self)</code>	Extract the bottom level from a dataset This extracts the bottom level from each netCDF file.
<code>DataSet.vertical_interp(self[, levels])</code>	Vertically interpolate a dataset based on given vertical levels This is calculated for each time step and grid cell
<code>DataSet.vertical_mean(self)</code>	Calculate the depth-averaged mean for each variable This is calculated for each time step and grid cell
<code>DataSet.vertical_min(self)</code>	Calculate the vertical minimum of variable values This is calculated for each time step and grid cell
<code>DataSet.vertical_max(self)</code>	Calculate the vertical maximum of variable values This is calculated for each time step and grid cell
<code>DataSet.vertical_range(self)</code>	Calculate the vertical range of variable values This is calculated for each time step and grid cell
<code>DataSet.vertical_sum(self)</code>	Calculate the vertical sum of variable values This is calculated for each time step and grid cell
<code>DataSet.vertical_cumsum(self)</code>	Calculate the vertical sum of variable values This is calculated for each time step and grid cell
<code>DataSet.invert_levels(self)</code>	Invert the levels of 3D variables This is calculated for each time step and grid cell
<code>DataSet.bottom_mask(self)</code>	Create a mask identifying the deepest cell without missing values.

14.9.1 nctoolkit.DataSet.surface

`DataSet.surface(self)`

Extract the top/surface level from a dataset This extracts the first vertical level from each file in a dataset.

Examples

If you wanted to extract the top vertical level of a dataset, do the following:

```
>>> data.surface()
```

This method is most useful for things like oceanic data, where this method will extract the sea surface.

14.9.2 nctoolkit.DataSet.bottom

`DataSet.bottom(self)`

Extract the bottom level from a dataset This extracts the bottom level from each netCDF file. Please note that for ensembles, it uses the first file to derive the index of the bottom level. Use `bottom_mask` for files when the bottom cell in netCDF files do not represent the actual bottom.

Examples

If you wanted to extract the bottom vertical level of a dataset, do the following:

```
>>> data.bottom()
```

This method is most useful for things like oceanic model data, where the bottom cell corresponds to the bottom of the ocean.

14.9.3 nctoolkit.DataSet.vertical_interp

`DataSet.vertical_interp(self, levels=None)`

Vertically interpolate a dataset based on given vertical levels This is calculated for each time step and grid cell

Parameters `levels` (*list, int or str*) – list of vertical levels, for example depths for an ocean model, to vertically interpolate to. These must be floats or ints.

Examples

If you wanted to vertically interpolate a dataset to 5 and 10 metres, you would do the following:

```
>>> data.vertical_interp([5,10])
```

This method is most useful for things like oceanic data, where you need to interpolate to certain depth levels. It will require that vertical levels are the same in every grid cell.

14.9.4 nctoolkit.DataSet.vertical_mean

`DataSet.vertical_mean(self)`

Calculate the depth-averaged mean for each variable This is calculated for each time step and grid cell

Examples

If you wanted to vertical mean of every variable in a dataset, you would do this:

```
>>> data.vertical_mean()
```

This method will calculate the vertical mean weighted by the thickness of each cell. Note that if cell thickness cannot be derived it will just average the values in each vertical cell.

14.9.5 nctoolkit.DataSet.vertical_min

`DataSet.vertical_min(self)`

Calculate the vertical minimum of variable values This is calculated for each time step and grid cell

Examples

If you wanted to vertical minimum of every variable in a dataset, you would do this:

```
>>> data.vertical_min()
```

14.9.6 nctoolkit.DataSet.vertical_max

`DataSet.vertical_max(self)`

Calculate the vertical maximum of variable values This is calculated for each time step and grid cell

Examples

If you wanted to vertical maximum of every variable in a dataset, you would do this:

```
>>> data.vertical_max()
```

14.9.7 nctoolkit.DataSet.vertical_range

`DataSet.vertical_range(self)`

Calculate the vertical range of variable values This is calculated for each time step and grid cell

Examples

If you wanted to range of values across all vertical levels of every variable in a dataset, you would do this:

```
>>> data.vertical_range()
```

14.9.8 nctoolkit.DataSet.vertical_sum

`DataSet.vertical_sum(self)`

Calculate the vertical sum of variable values This is calculated for each time step and grid cell

Examples

If you wanted to sum of values across all vertical levels of every variable in a dataset, you would do this:

```
>>> data.vertical_sum()
```

14.9.9 nctoolkit.DataSet.vertical_cumsum

`DataSet.vertical_cumsum(self)`

Calculate the vertical sum of variable values This is calculated for each time step and grid cell

Examples

If you wanted to calculate the cumulative sum of values across all vertical levels of every variable in a dataset, you would do this:

```
>>> data.vertical_sum()
```

The cumulative sum will be calculated from the first to the last vertical level. For example, in oceanic data it would start at the sea surface.

14.9.10 nctoolkit.DataSet.invert_levels

`DataSet.invert_levels(self)`

Invert the levels of 3D variables This is calculated for each time step and grid cell

Examples

If you wanted to invert the vertical levels, you would do this:

```
>>> data.invert_levels()
```


14.9.11 nctoolkit.DataSet.bottom_mask

`DataSet.bottom_mask(self)`

Create a mask identifying the deepest cell without missing values. This converts a dataset to a mask identifying which cell represents the bottom, for example the seabed. 1 identifies the deepest cell with non-missing values. Everything else is 0, or missing. At present this method only uses the first available variable from netCDF files, so it may not be suitable for all data

14.10 Rolling methods

<code>DataSet.rolling_mean(self[, window])</code>	Calculate a rolling mean based on a window
<code>DataSet.rolling_min(self[, window])</code>	Calculate a rolling minimum based on a window
<code>DataSet.rolling_max(self[, window])</code>	Calculate a rolling maximum based on a window
<code>DataSet.rolling_sum(self[, window])</code>	Calculate a rolling sum based on a window
<code>DataSet.rolling_range(self[, window])</code>	Calculate a rolling range based on a window

14.10.1 nctoolkit.DataSet.rolling_mean

`DataSet.rolling_mean(self, window=None)`

Calculate a rolling mean based on a window

Parameters = `int (window)` – The size of the window for the calculation of the rolling mean

Examples

If you wanted to calculate a rolling mean with the mean calculated over every 10 time steps, do the following:

```
>>> data.rolling_mean(10)
```

14.10.2 nctoolkit.DataSet.rolling_min

`DataSet.rolling_min(self, window=None)`

Calculate a rolling minimum based on a window

Parameters = `int (window)` – The size of the window for the calculation of the rolling minimum

Examples

If you wanted to calculate a rolling minimum with the minimum calculated over every 10 time steps, do the following:

```
>>> data.rolling_min(10)
```

14.10.3 nctoolkit.DataSet.rolling_max

`DataSet.rolling_max(self, window=None)`

Calculate a rolling maximum based on a window

Parameters = `int (window)` – The size of the window for the calculation of the rolling maximum

Examples

If you wanted to calculate a rolling maximum with the maximum calculated over every 10 time steps, do the following:

```
>>> data.rolling_max(10)
```

14.10.4 nctoolkit.DataSet.rolling_sum

`DataSet.rolling_sum(self, window=None)`

Calculate a rolling sum based on a window

Parameters = `int (window)` – The size of the window for the calculation of the rolling sum

Examples

If you wanted to calculate a rolling sum with the sum calculated over every 10 time steps, do the following:

```
>>> data.rolling_sum(10)
```

14.10.5 nctoolkit.DataSet.rolling_range

`DataSet.rolling_range(self, window=None)`

Calculate a rolling range based on a window

Parameters = `int (window)` – The size of the window for the calculation of the rolling range

Examples

If you wanted to calculate a rolling range with the range calculated over every 10 time steps, do the following:

```
>>> data.rolling_range(10)
```

14.11 Evaluation setting

`DataSet.run(self)`

Run all stored commands in a dataset

14.11.1 nctoolkit.DataSet.run

`DataSet.run(self)`

Run all stored commands in a dataset

Examples

If evaluation is lazy and you need to evaluate commands on a dataset, do the following:

```
>>> data.run()
```

14.12 Cleaning functions

14.13 Ensemble creation

`create_ensemble([path, recursive])`

Generate an ensemble

14.13.1 nctoolkit.create_ensemble

`nctoolkit.create_ensemble(path="", recursive=True)`

Generate an ensemble

Parameters

- **path** (*str*) – The directory to search for netCDF files
- **recursive** (*boolean*) – True/False depending on whether you want to search the path recursively. Defaults to True.

Returns A list of files

Return type list

Examples

If you wanted to recursively find all netCDF files available in a directory “data”, you would do this:

```
>>> import nctoolkit as nc
>>> nc.create_ensemble("data")
```

If you wanted to find the files in that directory and ignore subdirectories, you would instead do this:

```
>>> nc.create_ensemble("data", recursive = False)
```

14.14 Arithmetic methods

<code>DataSet.assign(self[, drop])</code>	Create new variables Existing columns that are re-assigned will be overwritten. :param drop: Set to True if you want existing variables to be removed once the new ones have been created. Defaults to False.
<code>DataSet.add(self[, x, var])</code>	Add to a dataset This will add a constant, another dataset or a netCDF file to the dataset. :param x: An int, float, single file dataset or netCDF file to add to the dataset. If a dataset or netCDF file is supplied, this must have only one variable, unless var is provided. The grids must be the same. :type x: int, float, DataSet or netCDF file :param var: A variable in the x to use for the operation :type var: str.
<code>DataSet.subtract(self[, x, var])</code>	Subtract from a dataset This will subtract a constant, another dataset or a netCDF file from the dataset. :param x: An int, float, single file dataset or netCDF file to subtract from the dataset. If a dataset or netCDF is supplied this must only have one variable, unless var is provided. The grids must be the same. :type x: int, float, DataSet or netCDF file :param var: A variable in the x to use for the operation :type var: str.
<code>DataSet.multiply(self[, x, var])</code>	Multiply a dataset This will multiply a dataset by a constant, another dataset or a netCDF file. :param x: An int, float, single file dataset or netCDF file to multiply the dataset by. If multiplying by a dataset or single file there must only be a single variable in it, unless var is supplied. The grids must be the same. :type x: int, float, DataSet or netCDF file :param var: A variable in the x to multiply the dataset by :type var: str.
<code>DataSet.divide(self[, x, var])</code>	Divide the data This will divide the dataset by a constant, another dataset or a netCDF file. :param x: An int, float, single file dataset or netCDF file to divide the dataset by. If a dataset or netCDF file is supplied, this must have only one variable, unless var is provided. The grids must be the same. :type x: int, float, DataSet or netCDF file :param var: A variable in the x to use for the operation :type var: str.

14.14.1 nctoolkit.DataSet.add

`DataSet.add(self, x=None, var=None)`

Add to a dataset This will add a constant, another dataset or a netCDF file to the dataset. :param x: An int, float, single file dataset or netCDF file to add to the dataset.

If a dataset or netCDF file is supplied, this must have only one variable, unless var is provided. The grids must be the same.

Parameters **var** (*str*) – A variable in the x to use for the operation

Examples

If you wanted to add 10 to all variables in a dataset, you would do the following:

```
>>> data.add(10)
```

To add the values in a dataset data2 from a dataset data1, you would do the following:

```
>>> data1.add(data2)
```

Grids in the datasets must match. Addition will occur in matching timesteps in data1 and data2. If there is only 1 timestep in data2, then the data from that timestep will be added to the data in all data1 time steps.

Adding the data from another netCDF file will work in the same way:

```
>>> data1.add("example.nc")
```

14.14.2 nctoolkit.DataSet.subtract

`DataSet.subtract` (*self*, *x=None*, *var=None*)

Subtract from a dataset This will subtract a constant, another dataset or a netCDF file from the dataset. :param x: An int, float, single file dataset or netCDF file to subtract from the dataset.

If a dataset or netCDF is supplied this must only have one variable, unless var is provided. The grids must be the same.

Parameters **var** (*str*) – A variable in the x to use for the operation

Examples

If you wanted to subtract 10 from all variables in a dataset, you would do the following:

```
>>> data.subtract(10)
```

To subtract the values in a dataset data2 from those in a dataset data1, you would do the following:

```
>>> data1.subtract(data2)
```

Grids in the datasets must match. Division will occur in matching timesteps in data1 and data2 if there are matching timesteps. If there is only 1 timestep in data2, then the data from that timestep in data2 will be subtracted from the data in all timesteps in data1.

Subtracting of the data from another netCDF file will work in the same way:

```
>>> data1.subtract("example.nc")
```

14.14.3 nctoolkit.DataSet.multiply

`DataSet.multiply(self, x=None, var=None)`

Multiply a dataset This will multiply a dataset by a constant, another dataset or a netCDF file. :param x: An int, float, single file dataset or netCDF file to multiply the dataset by.

If multiplying by a dataset or single file there must only be a single variable in it, unless var is supplied. The grids must be the same.

Parameters **var** (*str*) – A variable in the x to multiply the dataset by

Examples

If you wanted to multiply variables in a dataset by 10, you would do the following:

```
>>> data.multiply(10)
```

To multiply the values in a dataset by the values of variables in dataset data2, you would do the following:

```
>>> data1.multiply(data2)
```

Grids in the datasets must match. Multiplication will occur in matching timesteps in data1 and data2. If there is only 1 timestep in data2, then the data from that timestep in data2 will multiply the data in all timesteps in data1.

Multiplying a dataset by the data from another netCDF file will work in the same way:

```
>>> data1.multiply("example.nc")
```

14.14.4 nctoolkit.DataSet.divide

`DataSet.divide(self, x=None, var=None)`

Divide the data This will divide the dataset by a constant, another dataset or a netCDF file. :param x: An int, float, single file dataset or netCDF file to divide the dataset by.

If a dataset or netCDF file is supplied, this must have only one variable, unless var is provided. The grids must be the same.

Parameters **var** (*str*) – A variable in the x to use for the operation

Examples

If you wanted to dividie all variables in a dataset by 20, you would do the following:

```
>>> data.divide(10)
```

To divide values in a dataset by those in the dataset data2 from a dataset data1, you would do the following:

```
>>> data1.divide(data2)
```

Grids in the datasets must match. Division will occur in matching timesteps in data1 and data2. If there is only 1 timestep in data2, then the data from that timeste in data2 will divided the data in all data1 time steps.

Adding the data from another netCDF file will work in the same way:

```
>>> data1.divide("example.nc")
```

14.15 Ensemble statistics

<code>DataSet.ensemble_mean(self[, nco, ignore_time])</code>	Calculate an ensemble mean
<code>DataSet.ensemble_min(self[, nco, ignore_time])</code>	Calculate an ensemble min
<code>DataSet.ensemble_max(self[, nco, ignore_time])</code>	Calculate an ensemble maximum
<code>DataSet.ensemble_percentile(self[, p])</code>	Calculate an ensemble percentile This will calculate the percentiles for each time step in the files.
<code>DataSet.ensemble_range(self)</code>	Calculate an ensemble range The range is calculated for each time step; for example, if each file in the ensemble has 12 months of data the statistic will be calculated for each month.
<code>DataSet.ensemble_sum(self)</code>	Calculate an ensemble sum The sum is calculated for each time step; for example, if each file in the ensemble has 12 months of data the statistic will be calculated for each month.

14.15.1 nctoolkit.DataSet.ensemble_mean

`DataSet.ensemble_mean(self, nco=False, ignore_time=False)`

Calculate an ensemble mean

Parameters

- **nco** (*boolean*) – Do you want to use NCO for the calculation? Default is False, i.e. CDO is used. Modify default if run time is an issue.
- **ignore_time** (*boolean*) – If True the mean is calculated over all time steps. If False, the ensemble mean is calculated for each time steps; for example, if the ensemble is made up of monthly files the mean for each month will be calculated.

14.15.2 nctoolkit.DataSet.ensemble_min

`DataSet.ensemble_min(self, nco=False, ignore_time=False)`

Calculate an ensemble min

Parameters

- **nco** (*boolean*) – Do you want to use NCO for the calculation? Default is False, i.e. CDO is used. Modify default if run time is an issue.
- **ignore_time** (*boolean*) – If True the min is calculated over all time steps. If False, the ensemble min is calculated for each time steps; for example, if the ensemble is made up of monthly files the min for each month will be calculated.

14.15.3 nctoolkit.DataSet.ensemble_max

`DataSet.ensemble_max(self, nco=False, ignore_time=False)`

Calculate an ensemble maximum

Parameters

- **nco** (*boolean*) – Do you want to use NCO for the calculation? Default is False, i.e. CDO is used. Modify default if run time is an issue.
- **ignore_time** (*boolean*) – If True the max is calculated over all time steps. If False, the ensemble max is calculated for each time steps; for example, if the ensemble is made up of monthly files the max for each month will be calculated.

14.15.4 nctoolkit.DataSet.ensemble_percentile

`DataSet.ensemble_percentile(self, p=None)`

Calculate an ensemble percentile This will calculate the percentles for each time step in the files. For example, if you had an ensemble of files where each file included 12 months of data, it would calculate the percentile for each month.

Parameters **p** (*float or int*) – percentile to calculate. $0 \leq p \leq 100$.

14.15.5 nctoolkit.DataSet.ensemble_range

`DataSet.ensemble_range(self)`

Calculate an ensemble range The range is calculated for each time step; for example, if each file in the ensemble has 12 months of data the statistic will be calculated for each month.

14.15.6 nctoolkit.DataSet.ensemble_sum

`DataSet.ensemble_sum(self)`

Calculate an ensemble sum The sum is calculated for each time step; for example, if each file in the ensemble has 12 months of data the statistic will be calculated for each month.

14.16 Subsetting operations

<code>DataSet.crop(self[, lon, lat, nco, nco_vars])</code>	Crop to a rectangular longitude and latitude box
<code>DataSet.select(self, **kwargs)</code>	A method for subsetting datasets to specific variables, years, longitudes etc.
<code>DataSet.drop(self[, vars])</code>	Remove variables This will remove stated variables from files in the dataset.

14.16.1 nctoolkit.DataSet.crop

`DataSet.crop(self, lon=[-180, 180], lat=[-90, 90], nco=False, nco_vars=None)`

Crop to a rectangular longitude and latitude box

Parameters

- **lon** (*list*) – The longitude range to select. This must be two variables, between -180 and 180 when `nco = False`.
- **lat** (*list*) – The latitude range to select. This must be two variables, between -90 and 90 when `nco = False`.
- **nco** (*boolean*) – Do you want this to use NCO for cropping? Defaults to `False`, and uses CDO. Set to `True` if you want to call NCO. NCO is typically better at handling very large horizontal grids.
- **nco_vars** (*str or list*) – If using NCO, the variables you want to select

Examples

If you wanted to crop a dataset to longitudes between -40 and 30 and latitudes between -10 and 40, you would do the following:

```
>>> data.crop(lon = [-40, 30], lat = [-10, 40])
```

If you wanted to select only the northern hemisphere, the following will work:

```
>>> data.crop(lat = [0, 90])
```

14.16.2 nctoolkit.DataSet.select

`DataSet.select(self, *kwargs)`

A method for subsetting datasets to specific variables, years, longitudes etc. Operations are applied in the order supplied.

Parameters **kwargs* – Possible arguments: variables, years, months, seasons, timesteps, lon, lat

Note: this uses partial matches. So `year`, `month`, `var` etc. will also work

Each kwarg works as follows:

variables [str or list] A variable or list of variables to select

seasons [str] Seasons to select. One of “DJF”, “MAM”, “JJA”, “SON”.

months [list, range or int] Month(s) to select.

years [list, range or int] Years(s) to select. These should be integers

timesteps [list or int] time step(s) to select. For example, if you wanted the first time step set `times=0`.

Examples

If you want to select a single variable do the following:

```
>>> data.select(variable = "var")
```

If you want to select a list of variables, do this:

```
>>> data.select(variable = ["var1", "var2"])
```

If you want to select data for January, do the following:

```
>>> data.select(month = 1)
```

If you want to select a range of months, do the following:

```
>>> data.select(months = range(1, 7))
```

If you want to select a range of years, for example the 2010s, do the following:

```
>>> data.select(years = range(2010, 2020))
```

If you want to select the first two timesteps in a dataset, do the following:

```
>>> data.select(timesteps = [0, 1])
```

14.16.3 nctoolkit.DataSet.drop

`DataSet.drop(self, vars=None)`

Remove variables This will remove stated variables from files in the dataset.

Parameters **vars** (*str or list*) – Variable or variables to be removed from the dataset. Variables that are listed but not in the dataset will be ignored

Examples

If you wanted to remove a single variable ‘var1’ from a dataset data, you would do the following:

```
>>> data.drop('var')
```

If you wanted to remove a list of variables, you would do the following:

```
>>> data.drop(['var1', 'var2', 'var2'])
```

14.17 Time-based methods

`DataSet.set_date(self[, year, month, day, ...])`

Set the date in a dataset You should only do this if you have to fix/change a dataset with a single, not multiple dates.

`DataSet.shift(self, **kwargs)`

Shift method.

14.17.1 nctoolkit.DataSet.set_date

`DataSet.set_date(self, year=None, month=None, day=None, base_year=1900)`

Set the date in a dataset You should only do this if you have to fix/change a dataset with a single, not multiple dates.

Parameters

- **year** (*int*) – The year
- **month** (*int*) – The month
- **day** (*int*) – The day
- **base_year** (*int*) – The base year for time creation in the netCDF. Defaults to 1900.

14.17.2 nctoolkit.DataSet.shift

`DataSet.shift(self, **kwargs)`

Shift method. A wrapper for shift_days, shift_hours Operations are applied in the order supplied.

Parameters ***kwargs** – hours maps to shift_hours days maps to shift_days months maps to shift_months years maps to shift_years

Note: this uses partial matches. So hour, day, month, year will also work.

Examples

If you wanted to shift all times back 1 hour, you would do the following:

```
>>> data.shift(hours = -1)
```

If you wanted to shift all times forward 2 days, you would do the following:

```
>>> data.shift(days = 2)
```

If you wanted to shift all times forward 6 months, you would do the following:

```
>>> data.shift(months = 6)
```

If you wanted to shift all times forward 1 year, you would do the following:

```
>>> data.shift(years = 1)
```

This method will allow partial matches in arguments. So the following will do the same thing:

```
>>> data.shift(year = 2)
```

```
>>> data.shift(years = 2)
```

14.18 Interpolation and resampling methods

<code>DataSet.regrid(self[, grid, method, recycle])</code>	Regrid a dataset to a target grid
<code>DataSet.to_latlon(self[, lon, lat, res, ...])</code>	Regrid a dataset to a regular latlon grid
<code>DataSet.resample_grid(self[, factor])</code>	Resample the horizontal grid of a dataset
<code>DataSet.time_interp(self[, start, end, ...])</code>	Temporally interpolate variables based on date range and time resolution
<code>DataSet.timestep_interp(self[, steps])</code>	Temporally interpolate a dataset to given number of time steps between existing time steps

14.18.1 nctoolkit.DataSet.regrid

`DataSet.regrid(self, grid=None, method='bil', recycle=False)`

Regrid a dataset to a target grid

Parameters

- **grid** (*nctoolkit.DataSet*, *pandas data frame* or *netCDF file*) – The grid to remap to
- **method** (*str*) – Remapping method. Defaults to “bil”. Methods available are: bilinear - “bil”; nearest neighbour - “nn” - “nearest neighbour” bicubic interpolation - “bic” Distance-weighted average - “dis” First order conservative remapping - “con” Second order conservative remapping - “con2” Large area fraction remapping - “laf”

14.18.2 nctoolkit.DataSet.to_latlon

`DataSet.to_latlon(self, lon=None, lat=None, res=None, method='bil', recycle=False)`

Regrid a dataset to a regular latlon grid

Parameters

- **lon** (*list*) – 2 element list giving minimum and maximum longitude of target grid
- **lat** (*list*) – 2 element list giving minimum and maximum latitude of target grid
- **res** (*float*, *int* or *list*) – If float or int given, this will be the horizontal and vertical resolution of the target grid. If 2 element list is given, the first element is the longitudinal resolution and the second is the latitudinal resolution.
- **method** (*str*) – Remapping method. Defaults to “bil”. Methods available are: bilinear - “bil”; nearest neighbour - “nn” - “nearest neighbour” bicubic interpolation - “bic” Distance-weighted average - “dis” First order conservative remapping - “con” Second order conservative remapping - “con2” Large area fraction remapping - “laf”
- **recycle** (*bool*) – Do you want the grid and weights to be available for recycling and use in regrid? Defaults to False

14.18.3 nctoolkit.DataSet.resample_grid

`DataSet.resample_grid(self, factor=None)`

Resample the horizontal grid of a dataset

Parameters `factor` (*int*) – The resampling factor. Must be a positive integer. No interpolation occurs. Example: factor of 2 will sample every other grid cell

Examples

If you wanted to select every other grid cell, you could do the following:

```
>>> data.resample_grid(2)
```

14.18.4 nctoolkit.DataSet.time_interp

`DataSet.time_interp(self, start=None, end=None, resolution='monthly')`

Temporally interpolate variables based on date range and time resolution

Parameters

- **start** (*str*) – Start date for interpolation. Needs to be of the form YYYY/MM/DD or YYYY-MM-DD.
- **end** (*str*) – End date for interpolation. Needs to be of the form YYYY/MM/DD or YYYY-MM-DD. If end is not given interpolation will be to the final available time in the dataset.
- **resolution** (*str*) – Time steps used for interpolation. Needs to be “daily”, “weekly”, “monthly” or “yearly”. Defaults to monthly.

14.18.5 nctoolkit.DataSet.timestep_interp

`DataSet.timestep_interp(self, steps=None)`

Temporally interpolate a dataset to given number of time steps between existing time steps

Parameters `steps` (*int*) – Number of time steps to interpolate between existing time steps. For example, if you wanted to go from daily to hourly data you would set steps=24.

14.19 Masking methods

`DataSet.mask_box(self[, lon, lat])`

Mask a lon/lat box

14.19.1 nctoolkit.DataSet.mask_box

`DataSet.mask_box(self, lon=[- 180, 180], lat=[- 90, 90])`

Mask a lon/lat box

Parameters

- **lon** (*list*) – Longitude range to mask. Must be of the form: [lon_min, lon_max]
- **lat** (*list*) – Latitude range to mask. Must be of the form: [lat_min, lat_max]

14.20 Statistical methods

<code>DataSet.tmean(self[, over])</code>	Calculate the temporal mean of all variables
<code>DataSet.tmin(self[, over])</code>	Calculate the temporal minimum of all variables
<code>DataSet.tmedian(self[, over])</code>	Calculate the temporal median of all variables :param over: Time periods to average over.
<code>DataSet.tpercentile(self[, p, over])</code>	Calculate the temporal percentile of all variables
<code>DataSet.tmax(self[, over])</code>	Calculate the temporal maximum of all variables
<code>DataSet.tsum(self[, over])</code>	Calculate the temporal sum of all variables
<code>DataSet.trange(self[, over])</code>	Calculate the temporal range of all variables
<code>DataSet.tvariance(self[, over])</code>	Calculate the temporal variance of all variables
<code>DataSet.tstdev(self[, over])</code>	Calculate the temporal standard deviation of all variables
<code>DataSet.tcumsum(self)</code>	Calculate the temporal cumulative sum of all variables
<code>DataSet.cor_space(self[, var1, var2])</code>	Calculate the correlation correct between two variables in space This is calculated for each time step.
<code>DataSet.cor_time(self[, var1, var2])</code>	Calculate the correlation correct in time between two variables The correlation is calculated for each grid cell, ignoring missing values.
<code>DataSet.spatial_mean(self)</code>	Calculate the area weighted spatial mean for all variables This is performed for each time step.
<code>DataSet.spatial_min(self)</code>	Calculate the spatial minimum for all variables This is performed for each time step.
<code>DataSet.spatial_max(self)</code>	Calculate the spatial maximum for all variables This is performed for each time step.
<code>DataSet.spatial_percentile(self[, p])</code>	Calculate the spatial sum for all variables This is performed for each time step.
<code>DataSet.spatial_range(self)</code>	Calculate the spatial range for all variables This is performed for each time step.
<code>DataSet.spatial_sum(self[, by_area])</code>	Calculate the spatial sum for all variables This is performed for each time step.
<code>DataSet.centre(self[, by, by_area])</code>	Calculate the latitudinal or longitudinal centre for each year/month combination in files. This applies to each file in an ensemble. by : str Set to 'latitude' if you want the latitudinal centre calculated. 'longitude' for longitudinal. by_area : bool If the variable is a value/m2 type variable, set to True, otherwise set to False.
<code>DataSet.zonal_mean(self)</code>	Calculate the zonal mean for each year/month combination in files.

continues on next page

Table 21 – continued from previous page

<code>DataSet.zonal_min(self)</code>	Calculate the zonal minimum for each year/month combination in files.
<code>DataSet.zonal_max(self)</code>	Calculate the zonal maximum for each year/month combination in files.
<code>DataSet.zonal_range(self)</code>	Calculate the zonal range for each year/month combination in files.
<code>DataSet.meridional_mean(self)</code>	Calculate the meridional mean for each year/month combination in files.
<code>DataSet.meridional_min(self)</code>	Calculate the meridional minimum for each year/month combination in files.
<code>DataSet.meridional_max(self)</code>	Calculate the meridional maximum for each year/month combination in files.
<code>DataSet.meridional_range(self)</code>	Calculate the meridional range for each year/month combination in files.

14.20.1 nctoolkit.DataSet.tmean

`DataSet.tmean(self, over='time')`

Calculate the temporal mean of all variables

Parameters **over** (*str or list*) – Time periods to average over. Options are 'year', 'month', 'day'.

Examples

If you want to calculate mean over all time steps. Do the following:

```
>>> data.tmean()
```

If you want to calculate the mean for each year in a dataset, do this:

```
>>> data.tmean("year")
```

If you want to calculate the mean for each month in a dataset, do this:

```
>>> data.tmean("month")
```

If you want to calculate the mean for each month in each year in a dataset, do this:

```
>>> data.tmean(["year", "month"])
```

This method will also let you easily calculate climatologies. So, if you wanted to calculate a monthly climatological mean, you would do this:

```
>>> data.tmean("month")
```

A daily climatological mean would be the following:

```
>>> data.tmean("day")
```

14.20.2 nctoolkit.DataSet.tmin

`DataSet.tmin(self, over='time')`

Calculate the temporal minimum of all variables

Parameters **over** (*str or list*) – Time periods to average over. Options are ‘year’, ‘month’, ‘day’.

Examples

If you want to calculate minimum over all time steps. Do the following:

```
>>> data.tmin()
```

If you want to calculate the minimum for each year in a dataset, do this:

```
>>> data.tmin("year")
```

If you want to calculate the minimum for each month in a dataset, do this:

```
>>> data.tmin("month")
```

If you want to calculate the minimum for each month in each year in a dataset, do this:

```
>>> data.tmin(["year", "month"])
```

This method will also let you easily calculate climatologies. So, if you wanted to calculate a monthly climatological min, you would do this:

```
>>> data.tmin("month")
```

A daily climatological minimum would be the following:

```
>>> data.tmin("day")
```

14.20.3 nctoolkit.DataSet.tmedian

`DataSet.tmedian(self, over='time')`

Calculate the temporal median of all variables :param over: Time periods to average over. Options are ‘year’, ‘month’, ‘day’. :type over: str or list

Examples

If you want to calculate median over all time steps. Do the following:

```
>>> data.tmedian()
```

If you want to calculate the median for each year in a dataset, do this:

```
>>> data.tmedian("year")
```

If you want to calculate the median for each month in a dataset, do this:


```
>>> data.tmedian("month")
```

If you want to calculate the median for each month in each year in a dataset, do this:

```
>>> data.tmedian(["year", "month"])
```

This method will also let you easily calculate climatologies. So, if you wanted to calculate a monthly climatological median, you would do this:

```
>>> data.tmedian("month")
```

A daily climatological median would be the following:

```
>>> data.tmedian("day")
```

14.20.4 nctoolkit.DataSet.tpercentile

`DataSet.tpercentile(self, p=None, over='time')`

Calculate the temporal percentile of all variables

Parameters `p` (*float or int*) – Percentile to calculate

Examples

If you want to calculate the 20th percentile over all time steps. Do the following:

```
>>> data.tpercentile(20)
```

If you want to calculate the 20th percentile for each year in a dataset, do this:

```
>>> data.tpercentile(20)
```

14.20.5 nctoolkit.DataSet.tmax

`DataSet.tmax(self, over='time')`

Calculate the temporal maximum of all variables

Parameters `over` (*str or list*) – Time periods to average over. Options are 'year', 'month', 'day'.

Examples

If you want to calculate maximum over all time steps. Do the following:

```
>>> data.tmax()
```

If you want to calculate the maximum for each year in a dataset, do this:

```
>>> data.tmax("year")
```

If you want to calculate the maximum for each month in a dataset, do this:

```
>>> data.tmax("month")
```

If you want to calculate the maximum for each month in each year in a dataset, do this:

```
>>> data.tmax(["year", "month"])
```

This method will also let you easily calculate climatologies. So, if you wanted to calculate a monthly climatological max, you would do this:

```
>>> data.tmax("month")
```

A daily climatological maximum would be the following:

```
>>> data.tmax("day")
```

14.20.6 nctoolkit.DataSet.tsum

`DataSet.tsum(self, over='time')`

Calculate the temporal sum of all variables

14.20.7 nctoolkit.DataSet.trange

`DataSet.trange(self, over='time')`

Calculate the temporal range of all variables

Parameters **over** (*str or list*) – Time periods to average over. Options are ‘year’, ‘month’, ‘day’.

Examples

If you want to calculate range over all time steps. Do the following:

```
>>> data.trange()
```

If you want to calculate the range for each year in a dataset, do this:

```
>>> data.trange("year")
```

If you want to calculate the range for each month in a dataset, do this:

```
>>> data.trange("month")
```

If you want to calculate the range for each month in each year in a dataset, do this:

```
>>> data.trange(["year", "month"])
```

This method will also let you easily calculate climatologies. So, if you wanted to calculate a monthly climatological range, you would do this:

```
>>> data.trange("month")
```

A daily climatological range would be the following:

```
>>> data.trange( "day")
```

14.20.8 nctoolkit.DataSet.tvariance

`DataSet.tvariance` (*self*, *over*='time')

Calculate the temporal variance of all variables

Parameters **over** (*str* or *list*) – Time periods to average over. Options are ‘year’, ‘month’, ‘day’.

Examples

If you want to calculate variance over all time steps. Do the following:

```
>>> data.tvar()
```

If you want to calculate the variance for each year in a dataset, do this:

```
>>> data.tvar("year")
```

If you want to calculate the variance for each month in a dataset, do this:

```
>>> data.tvar("month")
```

If you want to calculate the variance for each month in each year in a dataset, do this:

```
>>> data.tvar(["year", "month"])
```

This method will also let you easily calculate climatologies. So, if you wanted to calculate a monthly climatological var, you would do this:

```
>>> data.tvar( "month")
```

A daily climatological variance would be the following:

```
>>> data.tvar( "day")
```

14.20.9 nctoolkit.DataSet.tstdev

`DataSet.tstdev` (*self*, *over*='time')

Calculate the temporal standard deviation of all variables

Parameters **over** (*str* or *list*) – Time periods to average over. Options are ‘year’, ‘month’, ‘day’.

Examples

If you want to calculate standard deviation over all time steps. Do the following:

```
>>> data.tstdev()
```

If you want to calculate the standard deviation for each year in a dataset, do this:

```
>>> data.tstdev("year")
```

If you want to calculate the standard deviation for each month in a dataset, do this:

```
>>> data.tstdev("month")
```

If you want to calculate the standard deviation for each month in each year in a dataset, do this:

```
>>> data.tstdev(["year", "month"])
```

This method will also let you easily calculate climatologies. So, if you wanted to calculate a monthly climatological var, you would do this:

```
>>> data.tstdev("month")
```

A daily climatological standard deviation would be the following:

```
>>> data.tstdev("day")
```

14.20.10 nctoolkit.DataSet.tcumsum

`DataSet.tcumsum(self)`

Calculate the temporal cumulative sum of all variables

Examples

If you want to calculate the cumulative sum for all variables over all timesteps, do this:

```
>>> data.tcumsum()
```

14.20.11 nctoolkit.DataSet.cor_space

`DataSet.cor_space(self, var1=None, var2=None)`

Calculate the correlation correct between two variables in space This is calculated for each time step. The correlation coefficient is calculated using values in all grid cells, ignoring missing values.

Parameters

- **var1** (*str*) – The first variable
- **var2** (*str*) – The second variable

Examples

If you wanted to calculate the spatial correlation coefficient between variables x and y in a dataset, you would do this:

```
>>> data.cor_space("x", "y")
```

The correlation coefficient will be calculated for each time step.

14.20.12 nctoolkit.DataSet.cor_time

`DataSet.cor_time(self, var1=None, var2=None)`

Calculate the correlation correct in time between two variables The correlation is calculated for each grid cell, ignoring missing values.

Parameters

- **var1** (*str*) – The first variable
- **var2** (*str*) – The second variable

Examples

If you wanted to calculate the temporal correlation coefficient between variables x and y in a dataset, you would do this:

```
>>> data.cor_space("x", "y")
```

The correlation coefficient will be calculated for each grid cell. This method will indicate how temporally correlated variables are in different spatial regions.

14.20.13 nctoolkit.DataSet.spatial_mean

`DataSet.spatial_mean(self)`

Calculate the area weighted spatial mean for all variables This is performed for each time step.

Examples

If you want to calculate the spatial mean for a dataset, just do the following:

```
>>> data.spatial_mean()
```

Note that this calculation will calculate the average using weights based on each cell's area. If cell areas cannot be calculated, it will take a straight average, and a warning will say this.

14.20.14 nctoolkit.DataSet.spatial_min

`DataSet.spatial_min(self)`

Calculate the spatial minimum for all variables This is performed for each time step.

Examples

If you want to calculate the spatial minimum for a dataset, just do the following:

```
>>> data.spatial_min()
```

14.20.15 nctoolkit.DataSet.spatial_max

`DataSet.spatial_max(self)`

Calculate the spatial maximum for all variables This is performed for each time step.

Examples

If you want to calculate the spatial maximum for a dataset, just do the following:

```
>>> data.spatial_max()
```

14.20.16 nctoolkit.DataSet.spatial_percentile

`DataSet.spatial_percentile(self, p=None)`

Calculate the spatial sum for all variables This is performed for each time step. :param p: Percentile to calculate. 0<=p<=100. :type p: int or float

Examples

If you want to calculate the median of each variable across space for a dataset, just do the following:

```
>>> data.spatial_percentile(50)
```

14.20.17 nctoolkit.DataSet.spatial_range

`DataSet.spatial_range(self)`

Calculate the spatial range for all variables This is performed for each time step.

Examples

If you want to calculate the range of each variable across space for a dataset, just do the following:

```
>>> data.spatial_max()
```

14.20.18 nctoolkit.DataSet.spatial_sum

`DataSet.spatial_sum(self, by_area=False)`

Calculate the spatial sum for all variables. This is performed for each time step.

Parameters `by_area` (*boolean*) – Set to True if you want to multiply the values by the grid cell area before summing over space. Default is False.

Examples

If you want to calculate the spatial sum each variable across space for a dataset, just do the following:

```
>>> data.spatial_sum()
```

By default, this method simply sums up each grid cell value. In some cases this is not suitable. For example, the values in each cell may be concentrations or values per square metre etc. In this case multiplying each cell value by the cell area is more suitable. Do the following:

```
>>> data.spatial_sum(by_area = True)
```

Each cell's value will be multiplied by the area of the cell (in square metres) prior to calculating the spatial sum.

14.20.19 nctoolkit.DataSet.centre

`DataSet.centre(self, by='latitude', by_area=False)`

Calculate the latitudinal or longitudinal centre for each year/month combination in files. This applies to each file in an ensemble. `by` : str

Set to 'latitude' if you want the latitudinal centre calculated. 'longitude' for longitudinal.

by_area [bool] If the variable is a value/m2 type variable, set to True, otherwise set to False.

14.20.20 nctoolkit.DataSet.zonal_mean

`DataSet.zonal_mean(self)`

Calculate the zonal mean for each year/month combination in files. This applies to each file in an ensemble.

Examples

If you want to calculate the zonal mean for a dataset, do the following:

```
>>> data.zonal_mean()
```

14.20.21 nctoolkit.DataSet.zonal_min

`DataSet.zonal_min(self)`

Calculate the zonal minimum for each year/month combination in files. This applies to each file in an ensemble.

Examples

If you want to calculate the zonal minimum for a dataset, do the following:

```
>>> data.zonal_min()
```

14.20.22 nctoolkit.DataSet.zonal_max

`DataSet.zonal_max(self)`

Calculate the zonal maximum for each year/month combination in files. This applies to each file in an ensemble.

Examples

If you want to calculate the zonal maximum for a dataset, do the following:

```
>>> data.zonal_max()
```

14.20.23 nctoolkit.DataSet.zonal_range

`DataSet.zonal_range(self)`

Calculate the zonal range for each year/month combination in files. This applies to each file in an ensemble.

Examples

If you want to calculate the zonal range for a dataset, do the following:

```
>>> data.zonal_range()
```


14.20.24 nctoolkit.DataSet.meridional_mean

`DataSet.meridional_mean(self)`

Calculate the meridional mean for each year/month combination in files. This applies to each file in an ensemble.

Examples

If you want to calculate the meridional mean for a dataset, do the following:

```
>>> data.meridional_mean()
```

14.20.25 nctoolkit.DataSet.meridional_min

`DataSet.meridional_min(self)`

Calculate the meridional minimum for each year/month combination in files. This applies to each file in an ensemble.

Examples

If you want to calculate the meridional minimum for a dataset, do the following:

```
>>> data.meridional_min()
```

14.20.26 nctoolkit.DataSet.meridional_max

`DataSet.meridional_max(self)`

Calculate the meridional maximum for each year/month combination in files. This applies to each file in an ensemble.

Examples

If you want to calculate the meridional maximum for a dataset, do the following:

```
>>> data.meridional_max()
```

14.20.27 nctoolkit.DataSet.meridional_range

`DataSet.meridional_range(self)`

Calculate the meridional range for each year/month combination in files. This applies to each file in an ensemble.

Examples

If you want to calculate the meridional range for a dataset, do the following:

```
>>> data.meridional_max()
```

14.21 Merging methods

<code>DataSet.merge(self[, match])</code>	Merge a multi-file ensemble into a single file Merging will occur based on the time steps in the first file.
<code>DataSet.merge_time(self)</code>	Time-based merging of a multi-file ensemble into a single file This method is ideal if you have the same data split over multiple files covering different data sets.

14.21.1 nctoolkit.DataSet.merge

`DataSet.merge(self, match=['year', 'month', 'day'])`

Merge a multi-file ensemble into a single file Merging will occur based on the time steps in the first file. This will only be effective if you want to merge files with the same times, but with different variables.

Parameters `match` (*list*, *str*) – a list or str stating what must match in the netCDF files. Defaults to year/month/day. This list must be some combination of year/month/day. An error will be thrown if the elements of time in match do not match across all netCDF files. The only exception is if there is a single date file in the ensemble.

14.21.2 nctoolkit.DataSet.merge_time

`DataSet.merge_time(self)`

Time-based merging of a multi-file ensemble into a single file This method is ideal if you have the same data split over multiple files covering different data sets.

14.22 Splitting methods

<code>DataSet.split(self[, by])</code>	Split the dataset Each file in the ensemble will be separated into new files based on the splitting argument.
--	---

14.22.1 nctoolkit.DataSet.split

`DataSet.split(self, by=None)`

Split the dataset Each file in the ensemble will be separated into new files based on the splitting argument.

Parameters `by` (*str*) – Available by arguments are 'year', 'month', 'yearmonth', 'season', 'day'. year will split files by year, month will split files by month, yearmonth will split files by year and month; season will split files by year, day will split files by day.

Examples

If you want to split each file into a dataset into a separate files for each year, do the following:

```
>>> data.split("year")
```

If you wanted to split by month, do the following:

```
>>> data.split("month")
```

14.23 Output and formatting methods

<code>DataSet.to_nc(self, out[, zip, overwrite])</code>	Save a dataset to a named file This will only work with single file datasets.
<code>DataSet.to_xarray(self[, decode_times, ...])</code>	Open a dataset as an xarray object
<code>DataSet.to_dataframe(self[, decode_times, ...])</code>	Open a dataset as a pandas data frame
<code>DataSet.zip(self)</code>	Zip the dataset This will compress the files within the dataset.
<code>DataSet.format(self[, ext])</code>	Zip the dataset This will compress the files within the dataset. This works lazily. :param ext: New format. Must be one of “nc”, “nc1”, “nc2”, “nc4” and “nc5”. netCDF = nc1 netCDF version 2 (64-bit offset) = nc2/netCDF4 (HDF5) = nc4 netCDF4-classic = nc4c netCDF version 5 (64-bit data) = nc5 :type ext: str.

14.23.1 nctoolkit.DataSet.to_nc

`DataSet.to_nc(self, out, zip=True, overwrite=False)`

Save a dataset to a named file This will only work with single file datasets.

Parameters

- **out** (*str*) – Output file name.
- **zip** (*boolean*) – True/False depending on whether you want to zip the file. Default is True.
- **overwrite** (*boolean*) – If out file exists, do you want to overwrite it? Default is False.

Examples

If you want to export a dataset to a netCDF file, do the following:

```
>>> data.to_nc("out.nc")
```

By default this file will be zipped. If you do not want it zipped, do this:

```
>>> data.to_nc("out.nc", zip = False)
```

By default this cannot overwrite files. If the output file exists, do the following:

```
>>> data.to_nc("out.nc", overwrite = True)
```

14.23.2 nctoolkit.DataSet.to_xarray

`DataSet.to_xarray(self, decode_times=True, cdo_times=False)`

Open a dataset as an xarray object

Parameters

- **decode_times** (*boolean*) – Set to False if you do not want xarray to decode the times. Default is True. If xarray cannot decode times, CDO will be used.
- **cdo_times** (*boolean*) – Set to True if you do not want CDO to decode the times

Returns to_xarray

Return type xarray.Dataset

Examples

If you want to convert a dataset to an xarray dataset, do the following:

```
>>> data.to_xarray()
```

This will return an xarray dataset.

If you do not want time to be decoded, do the following:

```
>>> data.to_xarray(decode_times = False)
```

14.23.3 nctoolkit.DataSet.to_dataframe

`DataSet.to_dataframe(self, decode_times=True, cdo_times=False)`

Open a dataset as a pandas data frame

Parameters

- **decode_times** (*boolean*) – Set to False if you do not want xarray to decode the times prior to conversion to data frame. Default is True.
- **cdo_times** (*boolean*) – Set to True if you do not want CDO to decode the times

Returns to_dataframe

Return type pandas.DataFrame

14.23.4 nctoolkit.DataSet.zip

`DataSet.zip(self)`

Zip the dataset This will compress the files within the dataset. This works lazily.

Examples

If you want to zip the files in a dataset, do the following:

```
>>> data.zip()
```

This will occur lazily, so will only occur after everything has been evaluated.

14.23.5 nctoolkit.DataSet.format

`DataSet.format(self, ext=None)`

Zip the dataset This will compress the files within the dataset. This works lazily. :param ext: New format. Must be one of “nc”, “nc1”, “nc2”, “nc4” and “nc5”.

netCDF = nc1 netCDF version 2 (64-bit offset) = nc2/nc netCDF4 (HDF5) = nc4 netCDF4-classi =
nc4c netCDF version 5 (64-bit data) = nc5

14.24 Miscellaneous methods

<code>DataSet.cell_area(self[, join])</code>	Calculate the area of grid cells.
<code>DataSet.first_above(self[, x])</code>	Identify the time step when a value is first above a threshold This will do the comparison with either a number, a Dataset or a netCDF file. :param x: An int, float, single file dataset or netCDF file to use for the threshold(s). If comparing with a dataset or single file there must only be a single variable in it. The grids must be the same. :type x: int, float, DataSet or netCDF file.
<code>DataSet.first_below(self[, x])</code>	Identify the time step when a value is first below a threshold This will do the comparison with either a number, a Dataset or a netCDF file. :param x: An int, float, single file dataset or netCDF file to use for the threshold(s). If comparing with a dataset or single file there must only be a single variable in it. The grids must be the same. :type x: int, float, DataSet or netCDF file.
<code>DataSet.last_above(self[, x])</code>	Identify the final time step when a value is above a threshold This will do the comparison with either a number, a Dataset or a netCDF file. :param x: An int, float, single file dataset or netCDF file to use for the threshold(s). If comparing with a dataset or single file there must only be a single variable in it. The grids must be the same. :type x: int, float, DataSet or netCDF file.

continues on next page

Table 25 – continued from previous page

<code>DataSet.last_above(self[, x])</code>	Identify the final time step when a value is above a threshold This will do the comparison with either a number, a Dataset or a netCDF file. :param x: An int, float, single file dataset or netCDF file to use for the threshold(s). If comparing with a dataset or single file there must only be a single variable in it. The grids must be the same. :type x: int, float, DataSet or netCDF file.
<code>DataSet.cdo_command(self[, command])</code>	Apply a cdo command
<code>DataSet.nco_command(self[, command, ensemble])</code>	Apply an nco command
<code>DataSet.compare_all(self[, expression])</code>	Compare all variables to a constant
<code>DataSet.gt(self, x)</code>	Method to calculate if variable in dataset is greater than that in another file or dataset This currently only works with single file datasets
<code>DataSet.lt(self, x)</code>	Method to calculate if variable in dataset is less than that in another file or dataset This currently only works with single file datasets
<code>DataSet.reduce_dims(self)</code>	Reduce dimensions of data This will remove any dimensions with only one value.
<code>DataSet.reduce_grid(self[, mask])</code>	Reduce the dataset to non-zero locations in a mask :param mask: single variable dataset or path to .nc file. The mask must have an identical grid to the dataset. :type mask: str or dataset.

14.24.1 nctoolkit.DataSet.cell_area

`DataSet.cell_area(self, join=True)`

Calculate the area of grid cells. Area of grid cells is given in square meters.

Parameters `join` (*boolean*) – Set to False if you only want the cell areas to be in the output. `join=True` adds the areas as a variable to the dataset. Defaults to True.

Examples

If you wanted to add the cell_areas as a new variable in a dataset, you would do the following:

```
>>> data.cell_area()
```

If you wanted to replace a dataset with the cell areas of that dataset, you would do the following:

```
>>> data.cell_area(join = False)
```

14.24.2 nctoolkit.DataSet.first_above

`DataSet.first_above(self, x=None)`

Identify the time step when a value is first above a threshold This will do the comparison with either a number, a Dataset or a netCDF file. :param x: An int, float, single file dataset or netCDF file to use for the threshold(s).

If comparing with a dataset or single file there must only be a single variable in it. The grids must be the same.

Examples

If you wanted to calculate the first time step where the value in a grid cell goes above 10, you would do the following

```
>>> data.first_above(10)
```

If you wanted to calculate the first time step where the value in a grid cell goes above that in another dataset, the following will work. Note that both datasets must have the same grid, and can only have single variables. The second dataset can, of course, only have one timestep.

```
>>> data.first_above(data1)
```

14.24.3 nctoolkit.DataSet.first_below

`DataSet.first_below(self, x=None)`

Identify the time step when a value is first below a threshold This will do the comparison with either a number, a Dataset or a netCDF file. :param x: An int, float, single file dataset or netCDF file to use for the threshold(s).

If comparing with a dataset or single file there must only be a single variable in it. The grids must be the same.

Examples

If you wanted to calculate the first time step where the value in a grid cell goes below 10, you would do the following

```
>>> data.first_below(10)
```

If you wanted to calculate the first time step where the value in a grid cell goes above that in another dataset, the following will work. Note that both datasets must have the same grid, and can only have single variables. The second dataset can, of course, only have one timestep.

```
>>> data.first_below(data1)
```

14.24.4 nctoolkit.DataSet.last_above

`DataSet.last_above(self, x=None)`

Identify the final time step when a value is above a threshold This will do the comparison with either a number, a Dataset or a netCDF file. :param x: An int, float, single file dataset or netCDF file to use for the threshold(s).

If comparing with a dataset or single file there must only be a single variable in it. The grids must be the same.

Examples

If you wanted to calculate the last time step where the value in a grid cell is above 10, you would do the following

```
>>> data.first_above(10)
```

If you wanted to calculate the last time step where the value in a grid cell goes above that in another dataset, the following will work. Note that both datasets must have the same grid, and can only have single variables. The second dataset can, of course, only have one timestep.

```
>>> data.first_above(data1)
```

14.24.5 nctoolkit.DataSet.cdo_command

`DataSet.cdo_command(self, command=None)`

Apply a cdo command

Parameters **command** (*string*) – cdo command to call. This command must be such that “cdo {command} infile outfile” will run.

14.24.6 nctoolkit.DataSet.nco_command

`DataSet.nco_command(self, command=None, ensemble=False)`

Apply an nco command

Parameters

- **command** (*string*) – nco command to call. This must be of a form such that “nco {command} infile outfile” will run.
- **ensemble** (*boolean*) – Set to True if you want the command to take all of the files as input. This is useful for ensemble methods.

14.24.7 nctoolkit.DataSet.compare_all

`DataSet.compare_all(self, expression=None)`

Compare all variables to a constant

Parameters `expression` (*str*) – This a regular comparison such as “<0”, “>0”, “==0”

Examples

If you wanted to identify grid cells with positive values you would do the following:

```
>>> data.compare_all(">0")
```

This will be calculated for each time step.

If you wanted to identify grid cells with negative values, you would do this

```
>>> data.compare_all("<0")
```

14.24.8 nctoolkit.DataSet.gt

`DataSet.gt(self, x)`

Method to calculate if variable in dataset is greater than that in another file or dataset This currently only works with single file datasets

Parameters `x` (*str or single file dataset*) – File path or nctoolkit dataset

14.24.9 nctoolkit.DataSet.lt

`DataSet.lt(self, x)`

Method to calculate if variable in dataset is less than that in another file or dataset This currently only works with single file datasets

Parameters `x` (*str or single file dataset*) – File path or nctoolkit dataset

14.24.10 nctoolkit.DataSet.reduce_dims

`DataSet.reduce_dims(self)`

Reduce dimensions of data This will remove any dimensions with only one value. For example, if only selecting one vertical level, the vertical dimension will be removed.

Examples

If you want to remove any dimensions that have only one value, do the following:

```
>>> data.reduce_dims("out.nc")
```

Note that this will work lazily. This method is most useful when you want to simplify datasets before exporting them to something like a pandas dataframe.

14.24.11 nctoolkit.DataSet.reduce_grid

`DataSet.reduce_grid(self, mask=None)`

Reduce the dataset to non-zero locations in a mask :param mask: single variable dataset or path to .nc file.

The mask must have an identical grid to the dataset.

14.25 Ecological methods

`DataSet.phenology(self[, var, metric, p])`

Calculate phenologies from a dataset Each file in an ensemble must only cover a single year, and ideally have all days.

14.25.1 nctoolkit.DataSet.phenology

`DataSet.phenology(self, var=None, metric=None, p=None)`

Calculate phenologies from a dataset Each file in an ensemble must only cover a single year, and ideally have all days. The method assumes datasets have daily resolution.

Parameters

- **var** (*str*) – Variable to analyze.
- **metric** (*str*) – Must be peak, middle, start or end. Peak is defined as the day of the maximum value. Middle is the day when the cumulative total of the variable first exceeds the cumulative total for the entire year. Start or end is defined as the first day when the cumulative total exceeds a percentile p of the maximum cumulative total.
- **p** (*str*) – Percentile to use for start or end.

PACKAGE INFO

This package was created by Robert Wilson at Plymouth Marine Laboratory (PML).

15.1 Acknowledgements

The current codebase of nctoolkit was developed using funding from the NERC Climate Linked Atlantic Sector Science programme ([NE/R015953/1](#)) and a combination of UK Research and Innovation (UKRI) and European Research Council (ERC) funded research projects.

15.2 Bugs and issues

If you identify bugs or issues with the package please raise an issue at PML's Marine Systems Modelling group's GitHub page [here](#) or contact nctoolkit's creator at rwi@pml.ac.uk.

15.3 Contributions welcome

The package is new, with new features being added each month. There remain a large number of features that could be added, especially for dealing with atmospheric data. If packages users are interested in contributing or suggesting new features they are welcome to raise and issue at the package's GitHub page or contact me.

PYTHON MODULE INDEX

n

`nctoolkit.append`, [45](#)

A

`add()` (*nctoolkit.DataSet* method), 56
`assign()` (*nctoolkit.DataSet* method), 47

B

`bottom()` (*nctoolkit.DataSet* method), 50
`bottom_mask()` (*nctoolkit.DataSet* method), 53

C

`cdo_command()` (*nctoolkit.DataSet* method), 84
`cell_area()` (*nctoolkit.DataSet* method), 82
`centre()` (*nctoolkit.DataSet* method), 75
`compare_all()` (*nctoolkit.DataSet* method), 85
`copy()` (*nctoolkit.DataSet* method), 44
`cor_space()` (*in module nctoolkit*), 44
`cor_space()` (*nctoolkit.DataSet* method), 72
`cor_time()` (*in module nctoolkit*), 44
`cor_time()` (*nctoolkit.DataSet* method), 73
`create_ensemble()` (*in module nctoolkit*), 55
`crop()` (*nctoolkit.DataSet* method), 61
`current()` (*nctoolkit.DataSet* property), 46

D

`divide()` (*nctoolkit.DataSet* method), 58
`drop()` (*nctoolkit.DataSet* method), 62

E

`ensemble_max()` (*nctoolkit.DataSet* method), 60
`ensemble_mean()` (*nctoolkit.DataSet* method), 59
`ensemble_min()` (*nctoolkit.DataSet* method), 59
`ensemble_percentile()` (*nctoolkit.DataSet* method), 60
`ensemble_range()` (*nctoolkit.DataSet* method), 60
`ensemble_sum()` (*nctoolkit.DataSet* method), 60

F

`first_above()` (*nctoolkit.DataSet* method), 83
`first_below()` (*nctoolkit.DataSet* method), 83
`format()` (*nctoolkit.DataSet* method), 81

G

`gt()` (*nctoolkit.DataSet* method), 85

H

`history()` (*nctoolkit.DataSet* property), 46

I

`invert_levels()` (*nctoolkit.DataSet* method), 52

L

`last_above()` (*nctoolkit.DataSet* method), 84
`levels()` (*nctoolkit.DataSet* property), 46
`lt()` (*nctoolkit.DataSet* method), 85

M

`mask_box()` (*nctoolkit.DataSet* method), 66
`merge()` (*in module nctoolkit*), 44
`merge()` (*nctoolkit.DataSet* method), 78
`merge_time()` (*nctoolkit.DataSet* method), 78
`meridional_max()` (*nctoolkit.DataSet* method), 77
`meridional_mean()` (*nctoolkit.DataSet* method), 77
`meridional_min()` (*nctoolkit.DataSet* method), 77
`meridional_range()` (*nctoolkit.DataSet* method), 77
`module`
 `nctoolkit.append`, 45
`months()` (*nctoolkit.DataSet* property), 46
`multiply()` (*nctoolkit.DataSet* method), 58

N

`nco_command()` (*nctoolkit.DataSet* method), 84
`nctoolkit.append`
 `module`, 45

O

`open_data()` (*in module nctoolkit*), 42
`open_thredds()` (*in module nctoolkit*), 43
`open_url()` (*in module nctoolkit*), 43
`options()` (*in module nctoolkit*), 41

P

`phenology()` (*nctoolkit.DataSet* method), 86
`plot()` (*nctoolkit.DataSet* method), 47

R

`reduce_dims()` (*nctoolkit.DataSet* method), 85
`reduce_grid()` (*nctoolkit.DataSet* method), 86
`regrid()` (*nctoolkit.DataSet* method), 64
`rename()` (*nctoolkit.DataSet* method), 48
`resample_grid()` (*nctoolkit.DataSet* method), 65
`rolling_max()` (*nctoolkit.DataSet* method), 54
`rolling_mean()` (*nctoolkit.DataSet* method), 53
`rolling_min()` (*nctoolkit.DataSet* method), 53
`rolling_range()` (*nctoolkit.DataSet* method), 54
`rolling_sum()` (*nctoolkit.DataSet* method), 54
`run()` (*nctoolkit.DataSet* method), 55

S

`select()` (*nctoolkit.DataSet* method), 61
`set_date()` (*nctoolkit.DataSet* method), 63
`set_longnames()` (*nctoolkit.DataSet* method), 49
`set_missing()` (*nctoolkit.DataSet* method), 48
`set_units()` (*nctoolkit.DataSet* method), 49
`shift()` (*nctoolkit.DataSet* method), 63
`size()` (*nctoolkit.DataSet* property), 46
`spatial_max()` (*nctoolkit.DataSet* method), 74
`spatial_mean()` (*nctoolkit.DataSet* method), 73
`spatial_min()` (*nctoolkit.DataSet* method), 74
`spatial_percentile()` (*nctoolkit.DataSet* method), 74
`spatial_range()` (*nctoolkit.DataSet* method), 74
`spatial_sum()` (*nctoolkit.DataSet* method), 75
`split()` (*nctoolkit.DataSet* method), 78
`start()` (*nctoolkit.DataSet* property), 46
`subtract()` (*nctoolkit.DataSet* method), 57
`sum_all()` (*nctoolkit.DataSet* method), 48
`surface()` (*nctoolkit.DataSet* method), 50

T

`tcumsum()` (*nctoolkit.DataSet* method), 72
`time_interp()` (*nctoolkit.DataSet* method), 65
`times()` (*nctoolkit.DataSet* property), 46
`timestep_interp()` (*nctoolkit.DataSet* method), 65
`tmax()` (*nctoolkit.DataSet* method), 69
`tmean()` (*nctoolkit.DataSet* method), 67
`tmedian()` (*nctoolkit.DataSet* method), 68
`tmin()` (*nctoolkit.DataSet* method), 68
`to_dataframe()` (*nctoolkit.DataSet* method), 80
`to_latlon()` (*nctoolkit.DataSet* method), 64
`to_nc()` (*nctoolkit.DataSet* method), 79
`to_xarray()` (*nctoolkit.DataSet* method), 80
`tpercentile()` (*nctoolkit.DataSet* method), 69
`trange()` (*nctoolkit.DataSet* method), 70
`tstddev()` (*nctoolkit.DataSet* method), 71
`tsum()` (*nctoolkit.DataSet* method), 70
`tvariance()` (*nctoolkit.DataSet* method), 71

V

`variables()` (*nctoolkit.DataSet* property), 45
`vertical_cumsum()` (*nctoolkit.DataSet* method), 52
`vertical_interp()` (*nctoolkit.DataSet* method), 50
`vertical_max()` (*nctoolkit.DataSet* method), 51
`vertical_mean()` (*nctoolkit.DataSet* method), 51
`vertical_min()` (*nctoolkit.DataSet* method), 51
`vertical_range()` (*nctoolkit.DataSet* method), 51
`vertical_sum()` (*nctoolkit.DataSet* method), 52

Y

`years()` (*nctoolkit.DataSet* property), 46

Z

`zip()` (*nctoolkit.DataSet* method), 81
`zonal_max()` (*nctoolkit.DataSet* method), 76
`zonal_mean()` (*nctoolkit.DataSet* method), 75
`zonal_min()` (*nctoolkit.DataSet* method), 76
`zonal_range()` (*nctoolkit.DataSet* method), 76