

---

# **wwdata Documentation**

***Release 0.1.0***

**Chaim De Mulder**

**Jun 12, 2018**



# CONTENTS

<b>1</b>	<b>wwdata</b>	<b>3</b>
1.1	Structure . . . . .	3
1.2	Examples . . . . .	5
1.3	Credits . . . . .	5
<b>2</b>	<b>Installation</b>	<b>7</b>
2.1	Stable release . . . . .	7
2.2	From sources . . . . .	7
<b>3</b>	<b>Usage</b>	<b>9</b>
<b>4</b>	<b>wwdata</b>	<b>11</b>
4.1	wwdata package . . . . .	11
4.1.1	Submodules . . . . .	11
4.1.2	wwdata.Class_HydroData module . . . . .	11
4.1.3	wwdata.Class_LabExperimBased module . . . . .	20
4.1.4	wwdata.Class_LabSensorBased module . . . . .	21
4.1.5	wwdata.Class_OnlineSensorBased module . . . . .	22
4.1.6	wwdata.data_reading_functions module . . . . .	28
4.1.7	wwdata.time_conversion_functions module . . . . .	30
4.1.8	Module contents . . . . .	31
<b>5</b>	<b>Contributing</b>	<b>33</b>
5.1	Types of Contributions . . . . .	33
5.1.1	Report Bugs . . . . .	33
5.1.2	Fix Bugs . . . . .	33
5.1.3	Implement Features . . . . .	33
5.1.4	Write Documentation . . . . .	33
5.1.5	Submit Feedback . . . . .	33
5.2	Get Started! . . . . .	34
5.3	Pull Request Guidelines . . . . .	34
5.4	Tips . . . . .	35
<b>6</b>	<b>Credits</b>	<b>37</b>
6.1	Development Lead . . . . .	37
6.2	Contributors . . . . .	37
6.3	Current funder . . . . .	37
<b>7</b>	<b>History</b>	<b>39</b>
7.1	0.1.0 (2017-09-01) . . . . .	39

<b>8 Indices and tables</b>	<b>41</b>
<b>Python Module Index</b>	<b>43</b>

Contents:



**WWDATA**

Data analysis package aimed at data obtained in the context of (waste)water

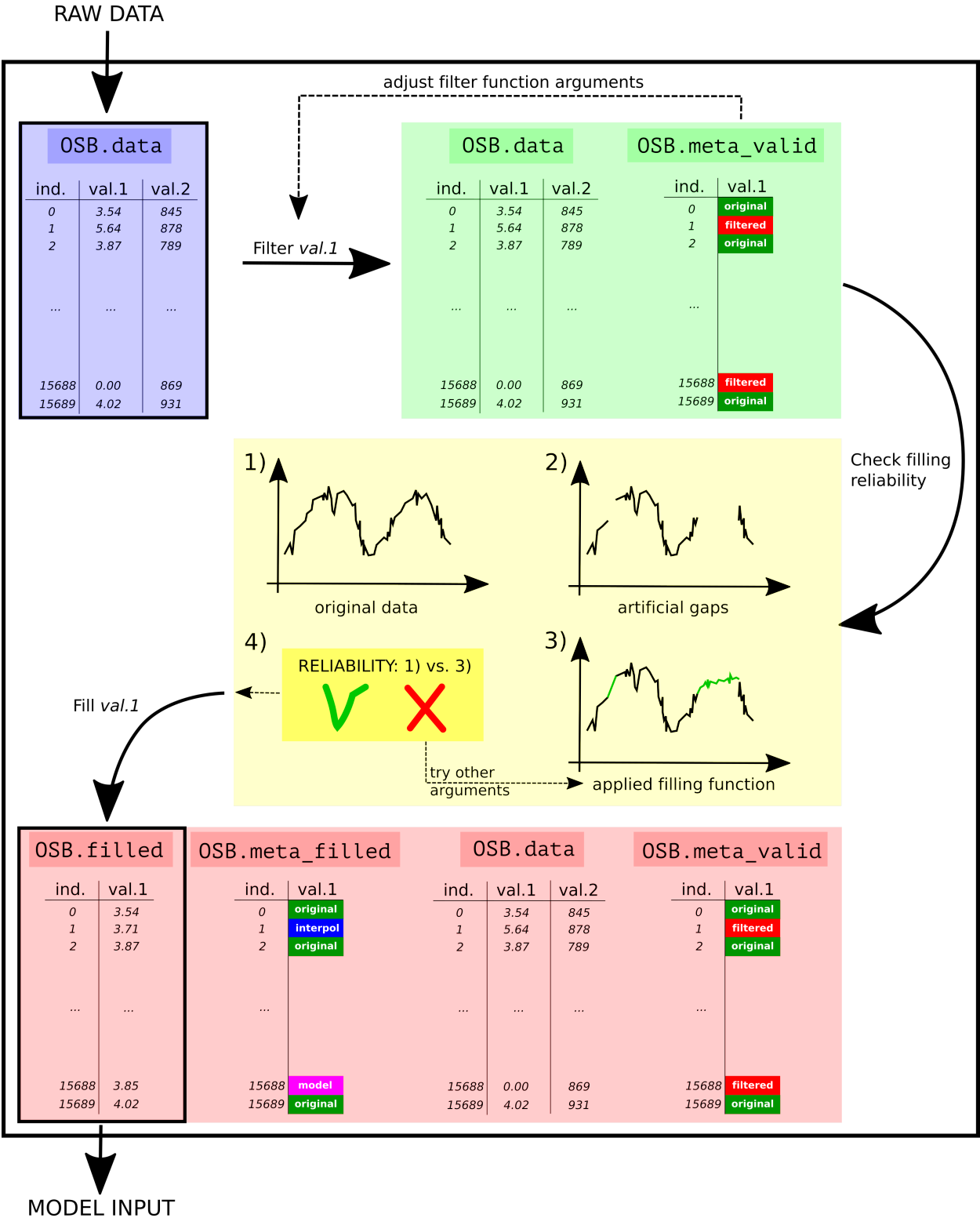
- Free software: GNU General Public License v3
- Documentation: <https://ugentbiomath.github.io/wwdata-docs/>
- Funding: Waterboard De Dommel
- Context: PhD research at BIOMATH, Ghent University

## 1.1 Structure

The package contains one class and three subclasses, all in separate .py files. Division in subclasses is based on the type of data:

- online data from full scale installations (OnlineSensorBased)
- online data from lab experiments (LabSensorBased)
- offline data obtained from lab experiments (LabExperimentBased).

Jupyter notebook files (.ipynb) illustrate the use of the available functions. The most developed class is the OnlineSensorBased one. The workflow of this class is shown in below Figure, where OSB represents an OnlineSensorBased object. Main premises are to never delete data but to tag it and to be able to check the reliability when gaps in datasets are filled.





## 1.2 Examples

For the workflow with code and more specific examples, check out the Showcase Jupyter Notebook(s) included as documentation of the package.

## 1.3 Credits

This package was created with support from [Cookiecutter](#) and the [audreyr/cookiecutter-pypackage](#) project template, as well as this [GitHub page](#), provided by [Daler](#) and explaining how to use sphinx documentation generation in combination with GitHub Pages.



## INSTALLATION

### 2.1 Stable release

To install wwdata, run this command in your terminal:

```
$ pip install wwdata
```

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

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

### 2.2 From sources

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

You can either clone the public repository:

```
$ git clone git://github.com/cdemulde/wwdata
```

Or download the [tarball](#):

```
$ curl -OL https://github.com/cdemulde/wwdata/tarball/master
```

Once you have a copy of the source, you can install it with:

```
$ python setup.py install
```



## USAGE

To use wwdata in a project:

```
import wwdata as ww
```

As the package was developed with the help of Jupyter Notebook (version 4.4.1) and visualisation is an important part of it, it is highly recommended to make use of a Jupyter Notebook when using the package.

Once the package is imported, the suggested workflow is as follows:

1. Read data and convert it to a pandas DataFrame (see [the pandas documentation](<https://pandas.pydata.org>) for more information); format it in the way you like.
2. Create an object of any of the three classes in the wwdata package, e.g. the OnlineSensorBased class:

```
data = ww.OnlineSensorBased(dataframe, timedata_column="time", data_type="WWTP",  
    ↪ experiment_tag="Data 2017")
```

3. Explore and format the data (convert to datetime, make absolute time index, make some plots...), e.g.:

```
data.to_datetime(time_column="time", time_format="%dd-%mm-%yy")  
data.get_avg()
```

4. Tag non-valid data points. The way to do this depends on the data you are working with, but the general approach would be:

```
data.tag_nan()  
data.moving_slope_filter(xdata="time", data_name="series1", cutoff=3, arange=['1/1/  
    ↪ 2017', '1/2/2017'])
```

i.e. to simply apply any of the filtering functions to the class object.

5. Apply any other functionalities to the object. In the below example, this is the filling of the gaps introduced by filtering data in the previous step (for details on the meaning of the arguments, please refer to the documentation provided within the source code):

```
data.fill_missing_interpolation("series1", range_=12, arange=['1/1/2017', '1/2/2017  
    ↪'])  
data.fill_missing_model("series1", "model_data_series", arange=['1/1/2017', '1/2/2017  
    ↪'])
```



## 4.1 wwdata package

### 4.1.1 Submodules

#### 4.1.2 wwdata.Class\_HydroData module

Class\_HydroData provides functionalities for handling data obtained in the context of (waste)water treatment.

Copyright (C) 2016 Chaim De Mulder

This program is free software: you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation, either version 3 of the License, or (at your option) any later version.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

You should have received a copy of the GNU General Public License along with this program. If not, see <http://www.gnu.org/licenses/>.

```
class wwdata.Class_HydroData.HydroData(data, timedata_column='index', data_type='WWTP',  
                                         experiment_tag='No tag given', time_unit=None,  
                                         units=[])
```

Bases: object

**timedata\_column**

*str* – name of the column containing the time data

**data\_type**

*str* – type of data provided

**experiment\_tag**

*str* – A tag identifying the experiment; can be a date or a code used by the producer/owner of the data.

**time\_unit**

*str* – The time unit in which the time data is given

**units**

*array* – The units of the variables in the columns

**absolute\_to\_relative** (*time\_data='index', unit='d', inplace=True, save\_abs=True, decimals=5*)

converts a pandas series with datetime timevalues to relative timevalues in the given unit, starting from 0

**Parameters**

- **time\_data** (*str*) – name of the column containing the time data. If this is the index column, just give ‘index’ (also default)
- **unit** (*str*) – unit to which to convert the time values (sec, min, hr or d)

#### Returns

- *None if inplace is True*
- *HydroData object if inplace is False*

#### **add\_to\_meta\_valid** (*column\_names*)

Adds (a) column(s) with the given column\_name(s) to the self.meta\_filled DataFrame, where all tags are set to ‘original’. This makes sure that also data that already is very reliable can be used further down the process (e.g. filling etc.)

**Parameters** **column\_names** (*array*) – array containing the names of the columns to add to the meta\_valied dataframe

#### **calc\_daily\_profile** (*column\_name, arange, quantile=0.9, plot=False, plot\_method='quantile', clear=False, only\_checked=False*)

Calculates a typical daily profile based on data from the indicated consecutive days. Also saves this average day, along with standard deviation and lower and upper percentiles as given in the arguments. Plotting is possible.

#### Parameters

- **column\_name** (*str*) – name of the column containing the data to calculate an average day for
- **arange** (*2-element array of ints*) – contains the beginning and end day of the period to use for average day calculation
- **quantile** (*float between 0 and 1*) – value to use for the calculation of the quantiles
- **plot** (*bool*) – plot or not
- **plot\_method** (*str*) – method to use for plotting. Available: “quantile” or “stdev”
- **clear** (*bool*) – wether or not to clear the key in the self.daily\_profile dictionary that is already present

**Returns** creates a dictionary self.daily\_profile containing information on the average day as calculated.

**Return type** None

#### **calc\_ratio** (*data\_1, data\_2, arange, only\_checked=False*)

Given two datasets or -columns, calculates the average ratio between the first and second dataset, within the given range. Also the standard deviation on this is calculated

#### Parameters

- **data\_1** (*str*) – name of the data column containing the data to be in the numerator of the ratio calculation
- **data\_2** (*str*) – name of the data column containing the data to be in the denominator of the ratio calculation
- **arange** (*array of two values*) – the range within which the ratio needs to be calculated
- **only\_checked** (*bool*) – if ‘True’, filtered values are excluded; default to ‘False’

**Returns**



- *The average ratio of the first data column over the second one within*
- *the given range and including the standard deviation*

**calc\_slopes** (*xdata*, *ydata*, *time\_unit=None*, *slope\_range=None*)

Calculates slopes for given xdata and data\_name; if a time unit is given as an argument, the time values (xdata) will first be converted to this unit, which will then be used to calculate the slopes with.

#### Parameters

- **xdata** (*str*) – name of the column containing the xdata for slope calculation (e.g. time). If 'index', the index is used as xdata. If datetime objects, a time\_unit is expected to calculate the slopes.
- **data\_name** (*str*) – name of the column containing the data\_name for slope calculation
- **time\_unit** (*str*) – time unit to be used for the slope calculation (in case this is based on time); if None, slopes are simply calculated based on the values given !! This value has no impact if the xdata column is the index and is not a datetime type. If that is the case, it is assumed that the user knows the unit of the xdata !!

**Returns** pandas Series object containing the slopes calculated for the chosen variable

**Return type** pd.Series

**compare\_ratio** (*data\_1*, *data\_2*, *arange*, *only\_checked=False*)

Compares the average ratios of two datasets in multiple different ranges and returns the most reliable one, based on the standard deviation on the ratio values

#### Parameters

- **data\_1** (*str*) – name of the data column containing the data to be in the numerator of the ratio calculation
- **data\_2** (*str*) – name of the data column containing the data to be in the denominator of the ratio calculation
- **arange** (*int*) – the range (in days) for which the ratios need to be calculated and compared
- **only\_checked** (*bool*) – if 'True', filtered values are excluded; default to 'False'

#### Returns

- *The average ratio within the range that has been found to be the most*
- *reliable one*

**drop\_index\_duplicates** ()

drop rows with a duplicate index. Also updates the meta\_valid dataframe

---

**Note:** It is assumed that the dropped rows contain the same data as their index- based duplicate, i.e. that no data is lost using the function.

---

**fill\_index** (*arange*, *index\_type='float'*)

function to fill in missing index values

**get\_avg** (*name=None*, *only\_checked=True*)

Gets the averages of all or certain columns in a dataframe

**Parameters** **name** (*array of str*) – name(s) of the column(s) containing the data to be averaged; defaults to ['none'] and will calculate average for every column

**Returns** pandas dataframe, containing the average slopes of all or certain columns

**Return type** `pd.DataFrame`

**get\_correlation** (*data\_1*, *data\_2*, *arange*, *zero\_intercept=False*, *only\_checked=False*, *plot=False*)

Calculates the linear regression coefficients that relate *data\_1* to *data\_2*

**Parameters**

- **and data\_2** (*data\_1*) – names of the data columns containing the data between which the correlation will be calculated.
- **arange** (*array*) – array containing the beginning and end value between which the correlation needs to be calculated
- **zero\_intercept** (*bool*) – indicates whether or not to assume a zero-intercept
- **only\_checked** (*bool*) – if ‘True’, filtered values are excluded from calculation and plotting; default to ‘False’ if a value in one column is filtered, the corresponding value in the second column also gets excluded!

**Returns**

- *the linear regression coefficients of the correlation, as well as the*
- *r-squared -value*

**get\_highs** (*data\_name*, *bound\_value*, *arange*, *method='percentile'*, *plot=False*)

creates a dataframe with tags indicating what indices have data-values higher than a certain value; example: the definition/tagging of rain events.

**Parameters**

- **data\_name** (*str*) – name of the column to execute the function on
- **bound\_value** (*float*) – the boundary value above which points will be tagged
- **arange** (*array of two values*) – the range within which high values need to be tagged
- **method** (*str (value or percentile)*) – when percentile, the bound value is a given percentile above which data points will be tagged, when value, bound\_values is used directly to tag data points.

**Returns**

**Return type** `None`

**get\_std** (*name=None*, *only\_checked=True*)

Gets the standard deviations of all or certain columns in a dataframe

**Parameters**

- **dataframe** (*pd.DataFrame*) – dataframe containing the columns to calculate the standard deviation for
- **name** (*array of str*) – name(s) of the column(s) containing the data to calculate standard deviation for; defaults to ['none'] and will calculate standard deviation for every column
- **plot** (*bool*) – if True, plots the calculated standard deviations, defaults to False

**Returns** pandas dataframe, containing the average slopes of all or certain columns

**Return type** `pd.DataFrame`

**head** (*n=5*)

piping pandas head function, see <https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.head.html> for documentation

**index** ()

piping pandas index function, see <http://pandas.pydata.org/pandas-docs/version/0.22/generated/pandas.Index.html> for documentation

**moving\_average\_filter** (*data\_name, window, cutoff\_frac, arange, clear=False, inplace=False, log\_file=None, plot=False, final=False*)

Filters out the peaks/outliers in a dataset by comparing its values to a smoothened representation of the dataset (Moving Average Filtering). The filtered values are replaced by NaN values.

#### Parameters

- **data\_name** (*str*) – name of the column containing the data that needs to be filtered
- **window** (*int*) – the number of values from the dataset that are used to take the average at the current point.
- **cutoff\_frac** (*float*) – the cutoff value (in fraction 0-1) to compare the data and smoothened data: a deviation higher than a certain percentage drops the data- point.
- **arange** (*array of two values*) – the range within which the moving average filter needs to be applied
- **clear** (*bool*) – if True, the tags added to datapoints before will be removed and put back to ‘original’.
- **inplace** (*bool*) – indicates whether a new dataframe is created and returned or whether the operations are executed on the existing dataframe (nothing is returned)
- **log\_file** (*str*) – string containing the directory to a log file to be written out when using this function
- **plot** (*bool*) – if true, a plot is made, comparing the original dataset with the new, filtered dataset
- **final** (*bool*) – if true, the values are actually replaced with nan values (either inplace or in a new hp object)

#### Returns

- *HydroData object (if inplace=False)* – the dataframe from which the double values of ‘data’ are removed
- *None (if inplace=True)*

**moving\_slope\_filter** (*xdata, data\_name, cutoff, arange, time\_unit=None, clear=False, inplace=False, log\_file=None, plot=False, final=False*)

Filters out datapoints based on the difference between the slope in one point and the next (sudden changes like noise get filtered out), based on a given cut off value. Replaces the dropped values with NaN values.

#### Parameters

- **xdata** (*str*) – name of the column containing the xdata for slope calculation (e.g. time). If ‘index’, the index is used as xdata. If datetime objects, a time\_unit is expected to calculate the slopes.
- **data\_name** (*str*) – name of the column containing the data that needs to be filtered
- **cutoff** (*int*) – the cutoff value to compare the slopes with to apply the filtering.
- **arange** (*array of two values*) – the range within which the moving slope filter needs to be applied

- **time\_unit** (*str*) – time unit to be used for the slope calculation (in case this is based on time); if None, slopes are calculated based on the values given
- **clear** (*bool*) – if True, the tags added to datapoints before will be removed and put back to ‘original’.
- **inplace** (*bool*) – indicates whether a new dataframe is created and returned or whether the operations are executed on the existing dataframe (nothing is returned)
- **log\_file** (*str*) – string containing the directory to a log file to be written out when using this function
- **plot** (*bool*) – if true, a plot is made, comparing the original dataset with the new, filtered dataset
- **final** (*bool*) – if true, the values are actually replaced with nan values (either inplace or in a new hp object)

#### Returns

- *HydroData object (if inplace=False)* – the dataframe from which the double values of ‘data’ are removed
- *None (if inplace=True)*
- *Creates*
- *\_\_\_\_\_*
- *A new column in the self.meta\_valid dataframe, containing a mask indicating*
- *what values are filtered*

**plot\_analysed** (*data\_name*, *time\_range*=‘default’, *only\_checked*=False)

plots the values and their types (original, filtered, filled) of a given column in the given time range.

#### Parameters

- **data\_name** (*str*) – name of the column containing the data to plot
- **time\_range** (*array of two values*) – the range within which the values are plotted; default is all
- **only\_checked** (*bool*) – if ‘True’, filtered values are excluded; default to ‘False’

#### Returns

**Return type** Plot

**replace** (*to\_replace*, *value*, *inplace*=False)

piping pandas replace function, see <http://pandas.pydata.org/pandas-docs/version/0.22/generated/pandas.DataFrame.replace.html> for documentation

**savgol** (*data\_name*, *window*=55, *polyorder*=2, *plot*=False, *inplace*=False)

Uses the scipy.signal Savitzky-Golay filter to smoothen the data of a column; The values are either replaced or a new dataframe is returned.

#### Parameters

- **data\_name** (*str*) – name of the column containing the data that needs to be filtered
- **window** (*int*) – the length of the filter window; default to 55
- **polyorder** (*int*) – The order of the polynomial used to fit the samples. polyorder must be less than window. default to 1

- **plot** (*bool*) – if true, a plot is made, comparing the original dataset with the new, filtered dataset
- **inplace** (*bool*) – indicates whether a new dataframe is created and returned or whether the operations are executed on the existing dataframe (nothing is returned)

**Returns**

- *HydroData object (if inplace=False)*
- *None (if inplace=True)*

**set\_index** (*keys*, *key\_is\_time=False*, *drop=True*, *inplace=False*, *verify\_integrity=False*, *save\_prev\_index=True*)

piping and extending pandas set\_index function, see [https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.set\\_index.html](https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.set_index.html) for documentation

**Notes**

**key\_is\_time** [bool] when true, the new index will be known as the time data from here on

(other arguments cfr pd.set\_index)

**Returns**

- *HydroData object (if inplace=False)*
- *None (if inplace=True)*

**set\_tag** (*tag*)

Sets the tag element of the HydroData object to the given tag

**Returns**

**Return type** None

**set\_time\_unit** (*unit*)

Sets the time\_unit element of the HydroData object to a given unit

**Returns**

**Return type** None

**set\_units** (*units*)

Set the units element of the HydroData object to a given dataframe

**simple\_moving\_average** (*arange*, *window*, *data\_name=None*, *inplace=False*, *plot=True*)

Calculate the Simple Moving Average of a dataserie from a dataframe, using a window within which the datavalues are averaged.

**Parameters**

- **arange** (*array of two values*) – the range within which the moving average needs to be calculated
- **window** (*int*) – the number of values from the dataset that are used to take the average at the current point. Defaults to 10
- **data\_name** (*str or array of str*) – name of the column(s) containing the data that needs to be smoothened. If None, smoothened data is computed for the whole dataframe. Defaults to None
- **inplace** (*bool*) – indicates whether a new dataframe is created and returned or whether the operations are executed on the existing dataframe (nothing is returned)

- **plot** (*bool*) – if True, a plot is given for comparison between original and smooth data

**Returns** either a new object (*inplace=False*) or an adjusted object, containing the smoothened data values

**Return type** *HydroData* (or subclass) object

**tag\_doubles** (*data\_name*, *bound*, *arange=None*, *clear=False*, *inplace=False*, *log\_file=None*, *plot=False*, *final=False*)

tags double values that subsequently occur in a measurement series. This is relevant in case a sensor has failed and produces a constant signal. A band is provided within which the signal can vary and still be filtered out

#### Parameters

- **data\_name** (*str*) – column name of the column from which double values will be sought
- **bound** (*float*) – boundary value of the band to use. When the difference between a point and the next one is smaller then the bound value, the latter datapoint is tagged as ‘filtered’.
- **arange** (*array of two values*) – the range within which double values need to be tagged
- **clear** (*bool*) – if True, the tags added to datapoints before will be removed and put back to ‘original’.
- **inplace** (*bool*) – indicates whether a new dataframe is created and returned or whether the operations are executed on the existing dataframe (nothing is returned). (This argument only comes into play when the ‘final’ argument is True)
- **log\_file** (*str*) – string containing the directory to a log file to be written out when using this function
- **plot** (*bool*) – whether or not to make a plot of the newly tagged data points
- **final** (*bool*) – if true, the values are actually replaced with nan values (either inplace or in a new hp object)

#### Returns

- *HydroData object (if inplace=False)* – the dataframe from which the double values of ‘data’ are removed or replaced
- *None (if inplace=True)*

**tag\_extremes** (*data\_name*, *arange=None*, *limit=0*, *method='below'*, *clear=False*, *plot=False*)

Tags values above or below a given limit.

#### Parameters

- **data\_name** (*str*) – name of the column containing the data to be tagged
- **arange** (*array of two values*) – the range within which extreme values need to be tagged
- **limit** (*int/float*) – limit below or above which values need to be tagged
- **method** (*'below' or 'above'*) – below tags all the values below the given limit, above tags the values above the limit
- **clear** (*bool*) – if True, the tags added before will be removed and put back to ‘original’.
- **plot** (*bool*) – whether or not to make a plot of the newly tagged data points

**Returns****Return type** None;**tag\_nan** (*data\_name*, *arange=None*, *clear=False*)

adds a tag 'filtered' in self.meta\_valid for every NaN value in the given column

**Parameters**

- **data\_name** (*str*) – column name of the column to apply the function to
- **arange** (*array of two values*) – the range within which nan values need to be tagged
- **clear** (*bool*) – when true, resets the tags in meta\_valid for the data in column *data\_name*

**Returns****Return type** None**tail** (*n=5*)piping pandas tail function, see <https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.tail.html> for documentation**to\_datetime** (*time\_column='index'*, *time\_format='%dd-%mm-%yy'*, *unit='D'*)

Piping and modifying pandas to\_datetime function

**Parameters**

- **time\_column** (*str*) – column name of the column where values need to be converted to date- time values. Default 'index' converts index values to datetime
- **time\_format** (*str*) – the format to use by to\_datetime function to convert strings to datetime format
- **unit** (*str*) – unit to use by to\_datetime function to convert int or float values to datetime format

**to\_float** (*columns='all'*)

convert values in given columns to float values

**Parameters** **columns** (*array of strings*) – column names of the columns where values need to be converted to floats**write** (*filename*, *filepath='/Users/chaimdemulder/Documents/Work/github/wwdata/docs'*, *method='all'*)**Parameters**

- **filepath** (*str*) – the path the output file should be saved to
- **filename** (*str*) – the name of the output file
- **method** (*str (all, filtered, filled)*) – depending on the method choice, different values will be written out: all values, only the filtered values or the filled values
- **for\_WEST** (*bool*) –
- **include\_units** (*bool*) –

**Returns****Return type** None; write an output filewwdata.Class\_HydroData.**total\_seconds** (*timedelta\_value*)

### 4.1.3 wwdata.Class\_LabExperimBased module

Class\_LabExperimBased provides functionalities for data handling of data obtained in lab experiments in the field of (waste)water treatment. Copyright (C) 2016 Chaim De Mulder

This program is free software: you can redistribute it and/or modify it under the terms of the GNU Affero General Public License as published by the Free Software Foundation, either version 3 of the License, or (at your option) any later version.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU Affero General Public License for more details.

You should have received a copy of the GNU Affero General Public License along with this program. If not, see <http://www.gnu.org/licenses/>.

```
class wwdata.Class_LabExperimBased.LabExperimBased(data,          timedata_column='index',
                                                    data_type='NAT',          experi-
                                                    ment_tag='No      tag      given',
                                                    time_unit=None)
```

Bases: `wwdata.Class_HydroData.HydroData`

Superclass for a HydroData object, expanding the functionalities with specific functions for data gathered in lab experiments.

**timedata\_column**

*str* – name of the column containing the time data

**data\_type**

*str* – type of the data provided

**experiment\_tag**

*str* – A tag identifying the experiment; can be a date or a code used by the producer/owner of the data.

**time\_unit**

*str* – The time unit in which the time data is given

**units**

*array* – The units of the variables in the columns

**add\_conc** (*column\_name*, *x*, *y*, *new\_name*='default')

calculates the concentration values of the given column and adds them as a new column to the DataFrame.

#### Parameters

- **column\_name** (*str*) – column with values
- **x** (*int*) – ...
- **y** (*int*) – ...
- **new\_name** (*str*) – name of the new column, default to 'column\_name + mg/L'

**calc\_slope** (*columns*, *time\_column*='h')

calculates the slope of the selected columns

#### Parameters

- **columns** (*array of strings*) – columns to calculate the slope for
- **time\_column** (*str*) – time used for calculation; default to 'h'

**check\_ph** (*ph\_column*='pH', *thresh*=0.4)

gives the maximal change in pH



**Parameters**

- **ph\_column** (*str*) – column with pH-values, default to 'pH'
- **threshold** (*int*) – threshold value for warning, default to '0.4'

**hours** (*time\_column='index'*)

calculates the hours from the relative values

**Parameters** **time\_column** (*string*) – column containing the relative time values; default to index

**in\_out** (*columns*)

(start\_values-end\_values)

**Parameters** **columns** (*array of strings*) –

**plot** (*columns, time\_column='index'*)

calculates the slope of the selected columns

**Parameters**

- **columns** (*array of strings*) – columns to plot
- **time\_column** (*str*) – time used for calculation; default to 'h'

**removal** (*columns*)

total removal of nitrogen (1-(end\_values/start\_values))

**Parameters** **columns** (*array of strings*) –

#### 4.1.4 wwdata.Class\_LabSensorBased module

Class\_LabSensorBased provides functionalities for data handling of data obtained in lab experiments with online sensors in the field of (waste)water treatment. Copyright (C) 2016 Chaim De Mulder

This program is free software: you can redistribute it and/or modify it under the terms of the GNU Affero General Public License as published by the Free Software Foundation, either version 3 of the License, or (at your option) any later version.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU Affero General Public License for more details.

You should have received a copy of the GNU Affero General Public License along with this program. If not, see <http://www.gnu.org/licenses/>.

**class** wwdata.Class\_LabSensorBased.**LabSensorBased** (*data, experiment\_tag='None'*)

Bases: *wwdata.Class\_HydroData.HydroData*

Superclass for a HydroData object, expanding the functionalities with specific functions for data gathered in lab experiments

**timedata\_column**

*str* – name of the column containing the time data

**data\_type**

*str* – type of data provided

**experiment\_tag**

*str* – A tag identifying the experiment; can be a date or a code used by the producer/owner of the data.

**time\_unit**

*str* – The time unit in which the time data is given

**units**

*array* – The units of the variables in the columns

**drop\_peaks** (*data\_name*, *cutoff*, *inplace=True*, *log\_file=None*)

Filters out the peaks larger than a cut-off value in a dataserries

**Parameters**

- **data\_name** (*str*) – the name of the column to use for the removal of peak values
- **cutoff** (*int*) – cut off value to use for the removing of peaks; values with an absolute value larger than this cut off will be removed from the data
- **inplace** (*bool*) – indicates whether a new dataframe is created and returned or whether the operations are executed on the existing dataframe (nothing is returned)
- **log\_file** (*str*) – string containing the directory to a log file to be written out when using this function

**Returns**

- *LabSensorBased object (if inplace=False)* – the dataframe from which the double values of ‘data’ are removed
- *None (if inplace=True)*

### 4.1.5 wwdata.Class\_OnlineSensorBased module

Class\_OnlineSensorBased provides functionalities for data handling of data obtained with online sensors in the field of (waste)water treatment. Copyright (C) 2016 Chaim De Mulder

This program is free software: you can redistribute it and/or modify it under the terms of the GNU Affero General Public License as published by the Free Software Foundation, either version 3 of the License, or (at your option) any later version.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU Affero General Public License for more details.

You should have received a copy of the GNU Affero General Public License along with this program. If not, see <http://www.gnu.org/licenses/>.

```
class wwdata.Class_OnlineSensorBased.OnlineSensorBased(data,                time-  
                                                         data_column='index',  
                                                         data_type='WWTP',    exper-  
                                                         iment_tag='No   tag   given',  
                                                         time_unit=None)
```

Bases: *wwdata.Class\_HydroData.HydroData*

Superclass for a HydroData object, expanding the functionalities with specific functions for data gathered at full scale by continous measurements

**timedata\_column**

*str* – name of the column containing the time data

**data\_type**

*str* – type of data provided

**experiment\_tag**

*str* – A tag identifying the experiment; can be a date or a code used by the producer/owner of the data.

**time\_unit**

*str* – The time unit in which the time data is given

**units**

*array* – The units of the variables in the columns

**add\_to\_filled** (*column\_names*)

*column\_names* : array

**calc\_daily\_average** (*column\_name, arange, plot=False*)

calculates the daily average of values in the given column and returns them as a 2D-array, containing the days and the average values on the respective days. Plotting is possible.

#### Parameters

- **column\_name** (*str*) – name of the column containing the data to calculate the average values for
- **arange** (*array of two values*) – the range within which daily averages need to be calculated
- **plot** (*bool*) – plot or not

**Returns** pandas dataframe, containing the daily means with standard deviations for the selected column

**Return type** pd.DataFrame

**calc\_total\_proportional** (*Q\_tot, Q, conc, new\_name='new', unit='mg/l', filled=False*)

Calculates the total concentration of an incoming flow, based on the given total flow and the separate incoming flows and concentrations

#### Parameters

- **Q\_tot** (*str*) – name of the column containing the total flow
- **Q** (*array of str*) – names of the columns containing the separate flows
- **conc** (*array of str*) – names of the columns containing the separate concentration values
- **new\_name** (*str*) – name of the column to be added
- **filled** (*bool*) – if true, use self.filled to calculate proportions from

---

**Note:** !!Order of columns in Q and conc must match!!

---

#### Returns

- *None*;
- *creates a hydropy object with added column for the proportional concentration*

**check\_filling\_error** (*nr\_iterations, data\_name, filling\_function, test\_data\_range, nr\_small\_gaps=0, max\_size\_small\_gaps=0, nr\_large\_gaps=0, max\_size\_large\_gaps=0, \*\*options*)

Uses the `_calculate_filling_error` function (refer to that docstring for more specific info) to calculate the error on the data points that are filled with a certain algorithm. Because `_calculate_filling_error` inserts random gaps, results differ every time it is used. `Check_filling_error` averages this out.

#### Parameters

- **nr\_iterations** (*int*) – The number of iterations to run for the calculation of the imputation error
- **data\_name** (*string*) – name of the column containing the data the filling reliability needs to be checked for.
- **function** (*filling*) – the name of the filling function to be tested for reliability
- **test\_data\_range** (*array of two values*) – an array containing the start and end point of the test data to be used. IMPORTANT: for testing filling with correlation, this range needs to include the range for correlation calculation and the filling range.
- **/ nr\_large\_gaps** (*nr\_small\_gaps*) – the number of small/large gaps to create in the dataset for testing
- **/ max\_size\_large\_gaps** (*max\_size\_small\_gaps*) – the maximum size of the gaps inserted in the data, expressed in data points
- **\*\*options** – Arguments for the filling function; refer to the relevant filling function to know what arguments to give

---

**Note:** When checking for the error on data filling, a period (arange argument) with mostly reliable data should be used. If for example large gaps are already present in the given data, this will heavily influence the returned error, as filled values will be compared with the values from the data gap.

---

**Returns** adds the average filling error the self.filling\_error dataframe

**Return type** None

**drop\_index\_duplicates** ()

drop rows with a duplicate index. Also updates the meta\_valid, meta\_filled and filled dataframes

---

**Note:** This operation assumes the dropped rows have the same data in them and therefor no data is lost.

---

**fill\_missing\_correlation** (*to\_fill, to\_use, arange, corr\_range, zero\_intercept=False, filtered\_only=True, plot=False, clear=False*)

Fills the missing values in a dataset (to\_fill), based on the correlation this data shows when comparing to other data (to\_use). This happens within the range given by arange.

**Parameters**

- **to\_fill** (*str*) – name of the column with data to fill
- **to\_use** (*str*) – name of the column to use, in combination with the given ratio, to fill in some of the missing data
- **arange** (*array of two values*) – the range within which missing/filtered values need to be replaced
- **corr\_range** (*array of two values*) – the range to use for the calculation of the correlation
- **filtered\_only** (*boolean*) – if True, fills only the datapoints labeled as filtered. If False, fills/replaces all datapoints in the given range
- **plot** (*bool*) – whether or not to plot the new dataset
- **clear** (*bool*) – whether or not to clear the previously filled values and start from the self.meta\_valid dataset again for this particular dataserie.

**Returns**

- *None*;
- *creates/updates self.filled, containing the adjusted dataset and updates*
- *meta\_filled with the correct labels.*

**fill\_missing\_daybefore** (*to\_fill*, *arange*, *range\_to\_replace*=[1, 4], *filtered\_only*=True, *plot*=False, *clear*=False)

Fills the missing values in a dataset (*to\_fill*), based on the data values from the day before the range starts. These data values are based on the *self.filled* dataset and therefor can contain filled datapoints as well. This happens within the range given by *arange*. **!! IMPORTANT !!** This function will not work on datasets with non-equidistant data points!

**Parameters**

- **to\_fill** (*str*) – name of the column with data to fill
- **arange** (*array of two values*) – the range within which missing/filtered values need to be replaced
- **range\_to\_replace** (*array of two int/float values*) – the minimum and maximum amount of time (i.e. min and max size of gaps in data) where missing datapoints can be replaced using this function, i.e. using values of the last day before measurements went bad.
- **filtered\_only** (*boolean*) – if True, fills only the datapoints labeled as filtered. If False, fills/replaces all datapoints in the given range
- **plot** (*bool*) – whether or not to plot the new dataset
- **clear** (*bool*) – whether or not to clear the previously filled values and start from the *self.meta\_valid* dataset again for this particular dataserie.

**Returns**

- *None*;
- *creates/updates self.filled, containing the adjusted dataset and updates*
- *meta\_filled with the correct labels.*

**fill\_missing\_interpolation** (*to\_fill*, *range\_*, *arange*, *method*='index', *plot*=False, *clear*=False)

Fills the missing values in a dataset (*to\_fill*), based specified interpolation algorithm (*method*). This happens only if the number of consecutive missing values is smaller than **range\_**.

**Parameters**

- **to\_fill** (*str*) – name of the column containing the data to be filled
- **range** (*int*) – the maximum range that the absence of values can be to still allow interpolation to fill in values
- **arange** (*array of two values*) – the range within which missing/filtered values need to be replaced
- **method** (*str*) – interpolation method to be used by the *.interpolate* function. See pandas docstrings for more info
- **plot** (*bool*) – whether or not to plot the new dataset
- **clear** (*bool*) – whether or not to clear the previously filled values and start from the *self.meta\_valid* dataset again for this particular dataserie.

**Returns**

- *None*;
- *creates/updates self.filled, containing the adjusted dataset and updates*
- *meta\_filled with the correct labels.*

**fill\_missing\_model** (*to\_fill*, *to\_use*, *arange*, *filtered\_only=True*, *unit='d'*, *plot=False*, *clear=False*)

Fills the missing values in a dataset (*to\_fill*), based on the modeled values given in *to\_use*. This happens within the range given by *arange*.

**Parameters**

- **to\_fill** (*str*) – name of the column with data to fill
- **to\_use** (*pd.Series*) – pandas series containing the modeled data with which the filtered data can be replaced
- **arange** (*array of two values*) – the range within which missing/filtered values need to be replaced
- **filtered\_only** (*boolean*) – if True, fills only the datapoints labeled as filtered. If False, fills/replaces all datapoints in the given range
- **unit** (*str*) – the unit in which the modeled values are given; datetime values will be converted to values with that unit. Possible: sec, min, hr, d
- **plot** (*bool*) – whether or not to plot the new dataset
- **clear** (*bool*) – whether or not to clear the previously filled values and start from the self.meta\_valid dataset again for this particular dataserie.

**Returns**

- *None*;
- *creates/updates self.filled, containing the adjusted dataset and updates*
- *meta\_filled with the correct labels.*

**fill\_missing\_ratio** (*to\_fill*, *to\_use*, *ratio*, *arange*, *filtered\_only=True*, *plot=False*, *clear=False*)

Fills the missing values in a dataset (*to\_fill*), based on the ratio this data shows when comparing to other data (*to\_use*). This happens within the range given by *arange*.

**Parameters**

- **to\_fill** (*str*) – name of the column with data to fill
- **to\_use** (*str*) – name of the column to use, in combination with the given ratio, to fill in some of the missing data
- **ratio** (*float*) – ratio to multiply the *to\_use* data with to obtain data for filling in the *to\_fill* data column
- **arange** (*array of two values*) – the range within which missing/filtered values need to be replaced
- **filtered\_only** (*boolean*) – if True, fills only the datapoints labeled as filtered. If False, fills/replaces all datapoints in the given range
- **plot** (*bool*) – whether or not to plot the new dataset
- **clear** (*bool*) – whether or not to clear the previously filled values and start from the self.meta\_valid dataset again for this particular dataserie.

**Returns**

- *None*;
- *creates/updates self.filled, containing the adjusted dataset and updates*
- *meta\_filled with the correct labels.*

**fill\_missing\_standard** (*to\_fill, arange, filtered\_only=True, plot=False, clear=False*)

Fills the missing values in a dataset (*to\_fill*), based on the average daily profile calculated by `calc_daily_profile()`. This happens within the range given by *arange*.

**Parameters**

- **to\_fill** (*str*) – name of the column with data to fill
- **arange** (*array of two values*) – the range within which missing/filtered values need to be replaced
- **filtered\_only** (*boolean*) – if True, fills only the datapoints labeled as filtered. If False, fills/replaces all datapoints in the given range
- **plot** (*bool*) – whether or not to plot the new dataset
- **clear** (*bool*) – whether or not to clear the previously filled values and start from the `self.meta_valid` dataset again for this particular dataseries.

**Returns**

- *None*;
- *creates/updates self.filled, containing the adjusted dataset and updates*
- *meta\_filled with the correct labels.*

`wwdata.Class_OnlineSensorBased.absolute_to_relative` (*series, start\_date, unit='d', decimals=5*)

converts a pandas series with datetime timevalues to relative timevalues in the given unit, starting from *start\_date*

**Parameters**

- **series** (*pd.Series*) – series of datetime of comparable values
- **unit** (*str*) – unit to which to convert the time values (sec, min, hr or d)
- **output** –
- -----

`wwdata.Class_OnlineSensorBased.drop_peaks` (*self, data\_name, cutoff, inplace=True, log\_file=None*)

Filters out the peaks larger than a cut-off value in a dataseries

**Parameters**

- **data\_name** (*str*) – the name of the column to use for the removal of peak values
- **cutoff** (*int*) – cut off value to use for the removing of peaks; values with an absolute value larger than this cut off will be removed from the data
- **inplace** (*bool*) – indicates whether a new dataframe is created and returned or whether the operations are executed on the existing dataframe (nothing is returned)
- **log\_file** (*str*) – string containing the directory to a log file to be written out when using this function

**Returns**

- *LabSensorBased object (if inplace=False)* – the dataframe from which the double values of ‘data’ are removed
- *None (if inplace=True)*

`wwdata.Class_OnlineSensorBased.find_nearest_time(value, df, column)`

Returns the (time) value in a dataframe column nearest to a given value

#### Parameters

- **value** (*float*) – time value to find the closest value for in ‘df’
- **df** (*pd.DataFrame*) – dataframe to use
- **column** (*str*) – column to check ‘value’ against

`wwdata.Class_OnlineSensorBased.go_WEST(raw_data, time_data, WEST_name_conversion)`

Saves a WEST compatible file (influent or other inputs)

#### Parameters

- **raw\_data** (*str or pd DataFrame*) –
- **time\_data** –
- **WEST\_name\_conversion** (*pd DataFrame with column names: WEST, units and RAW*) – dataframe containing three columns: the column names for the WEST-compatible file, the units to appear in the WEST-compatible file and the column names of the raw data file.
- **output** –
- **-----** –
- **None** –

`wwdata.Class_OnlineSensorBased.total_seconds(timedelta_value)`

`wwdata.Class_OnlineSensorBased.vlookup_day(value, df, column)`

Returns the dataframe index of a given value

## 4.1.6 wwdata.data\_reading\_functions module

`data_reading_functions` provides functionalities for data reading in the context of the `wwdata` package. Copyright (C) 2016 Chaim De Mulder

This program is free software: you can redistribute it and/or modify it under the terms of the GNU Affero General Public License as published by the Free Software Foundation, either version 3 of the License, or (at your option) any later version.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU Affero General Public License for more details.

You should have received a copy of the GNU Affero General Public License along with this program. If not, see <http://www.gnu.org/licenses/>.

`wwdata.data_reading_functions.find_and_replace(path, ext, replace)`

Finds the files with a certain extension in a directory and applies a find- replace action to those files. Removes the old files and produces files with a prefix stating the replacing value.

#### Parameters

- **path** (*str*) – the path name of the directory to apply the function to



- **ext** (*str*) – the extension of the files to be searched (excel, text or csv)
- **replace** (*array of str*) – the first value of replace is the string to be replaced by the second value of replace.

`wwdata.data_reading_functions.join_files` (*path, files, ext='text', sep=',', comment='#', encoding='utf8', decimal='.'*)

Reads all files in a given directory, joins them and returns one `pd.dataframe`

#### Parameters

- **path** (*str*) – path to the folder that contains the files to be joined
- **files** (*list*) – list of files to be joined, must be the same extension
- **ext** (*str*) – extension of the files to read; possible: excel, text, csv
- **sep** (*str*) – the separating element (e.g. , or ) necessary when reading csv-files
- **comment** (*str*) – comment symbol used in the files
- **sort** (*array of bool and str*) – if first element is true, apply the sort function to sort the data based on the tags in the column mentioned in the second element of the sort array

**Returns** pandas dataframe containin concatenated files in the given directory

**Return type** `pd.dataframe`

`wwdata.data_reading_functions.list_files` (*path, ext*)

Returns a list of files in a certain folder ('path') with a certain extension ('ext')

#### Parameters

- **path** (*str*) – path to the folder containing the files to be listed
- **ext** (*str*) – extension of the files to be listed; current options are 'excel','text' or 'csv'

`wwdata.data_reading_functions.read_mat` (*path*)

TO DO Reads in .mat datafiles and returns them as `pd.DataFrame` <http://stackoverflow.com/questions/24762122/read-matlab-data-file-into-python-need-to-export-to-csv>

`wwdata.data_reading_functions.remove_empty_lines` (*path, ext*)

Removes the empty lines from files in a certain folder ('path') and with a certain extension ('ext')

#### Parameters

- **path** (*str*) – path to the folder containing the files in which empty lines need to be removed
- **ext** (*str*) – extension of the files in which empty lines need to be removed; current options are 'excel','text' or 'csv'

`wwdata.data_reading_functions.sort_data` (*data, based\_on, reset\_index=[False, 'new\_index\_name'], convert\_to\_timestamp=[True, 'time\_name', '%d.%m.%Y %H:%M:%S']*)

Sorts a dataset based on values in one of the columns and splits them in different dataframes, returned in the form of one dictionary

#### Parameters

- **data** (*pd.dataframe*) – the dataframe containing the data that needs to be sorted
- **based\_on** (*str*) – the name of the column that contains the names or values the sorting should be based on

- **reset\_index** (*[bool, str]*) – array indicating if the index of the sorted datasets should be reset to a new one; if first element is true, the second element is the title of the column to use as new index; default: False

**Returns** A dictionary of pandas dataframes with as labels those acquired from the `based_on` column

**Return type** dict

```
wwdata.data_reading_functions.write_to_WEST(df, file_normal, file_west, units,
                                           filepath='/Users/chaimdemulder/Documents/Work/github/wwdata/do',
                                           fillna=True)
```

writes a text-file that is compatible with WEST. Adds the units as they are given in the 'units' argument.

**Parameters**

- **df** (*pd.DataFrame*) – the dataframe to write to WEST
- **file\_normal** (*str*) – name of the original file to write, not yet compatible with WEST
- **file\_west** (*str*) – name of the file that needs to be WEST compatible
- **units** (*array of strings*) – array containing the units for the respective columns in df
- **filepath** (*str*) – directory to save the files in; defaults to the current one
- **fillna** (*bool*) – when True, replaces nan values with 0 values (this might avoid WEST problems later one).

**Returns**

**Return type** None; writes files

## 4.1.7 wwdata.time\_conversion\_functions module

`time_conversion_functions` provides functionalities for converting certain types of time data to other types, in the context of the `wwdata` package. Copyright (C) 2016 Chaim De Mulder

This program is free software: you can redistribute it and/or modify it under the terms of the GNU Affero General Public License as published by the Free Software Foundation, either version 3 of the License, or (at your option) any later version.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU Affero General Public License for more details.

You should have received a copy of the GNU Affero General Public License along with this program. If not, see <http://www.gnu.org/licenses/>.

```
wwdata.time_conversion_functions.get_absolute_time(value,
                                                  date_type='WindowsDateSystem')
```

Converts a coded time to the absolute date at which the experiment was conducted.

**Parameters**

- **value** (*int or float*) – a code for a certain point in time
- **date\_type** (*str*) – type of coding used for the time point, probably depending on the software which was used to acquire the data, e.g. the Windows Date System (here as default, see also: <http://excelsemipro.com/2010/08/date-and-time-calculation-in-excel/>)

**Returns** python datetime object

**Return type** datetime.datetime

```
wwdata.time_conversion_functions.make_datetime(array)
```

**Parameters with elements** (*array*) – 0: year (yy) 1: month (mm) 2: day in month (dd) 3: hour (h or hh) 4: minutes (minmin)

```
wwdata.time_conversion_functions.make_month_day_array()
```

Returns a dataframe containing two columns, one with the number of the month, one with the day of the month. Useful in creating datetime objects based on e.g. date serial numbers from the Window Date System (<http://excelsemipro.com/2010/08/date-and-time-calculation-in-excel/>)

**Returns** dataframe with number of the month and number of the day of the month for a whole year

**Return type** pd.DataFrame

```
wwdata.time_conversion_functions.timedelta_to_abs(timedelta, unit='d')
```

*timedelta* : array of timedelta values

```
wwdata.time_conversion_functions.to_datetime_singlevalue(time)
```

In case *timedata* is in a string format, to convert it to a datetime object, it needs to be in the right format, e.g. dd-mm-yyyy hh:mm:ss (so two of each) This function takes care of that, to a certain extent.

```
wwdata.time_conversion_functions.to_days(timedelta)
```

*timedelta* : timedelta value

```
wwdata.time_conversion_functions.to_hours(timedelta)
```

*timedelta* : timedelta value

```
wwdata.time_conversion_functions.to_minutes(timedelta)
```

*timedelta* : timedelta value

```
wwdata.time_conversion_functions.to_seconds(timedelta)
```

*timedelta* : timedelta value

## 4.1.8 Module contents

Top-level package for wwdata.



## CONTRIBUTING

Contributions are welcome, and they are greatly appreciated! Every little bit helps, and credit will always be given. You can contribute in many ways:

### 5.1 Types of Contributions

#### 5.1.1 Report Bugs

Report bugs at <https://github.com/cdemulde/wwdata/issues>.

If you are reporting a bug, please include:

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

#### 5.1.2 Fix Bugs

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

#### 5.1.3 Implement Features

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

#### 5.1.4 Write Documentation

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

#### 5.1.5 Submit Feedback

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

If you are proposing a feature:

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

## 5.2 Get Started!

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

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

```
$ git clone git@github.com:your_name_here/wwdata.git
```

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

```
$ mkvirtualenv wwdata
$ cd wwdata/
$ python setup.py develop
```

4. Create a branch for local development:

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

Now you can make your changes locally.

5. When you're done making changes, check that your changes pass *flake8* and the tests, including testing other Python versions with *tox*:

```
$ flake8 wwdata tests
$ python setup.py test or py.test
$ tox
```

To get *flake8* and *tox*, just *pip* install them into your virtualenv.

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

```
$ git add .
$ git commit -m "Your detailed description of your changes."
$ git push origin name-of-your-bugfix-or-feature
```

7. Submit a pull request through the GitHub website.

## 5.3 Pull Request Guidelines

Before you submit a pull request, check that it meets these guidelines:

1. The pull request should include tests.
2. If the pull request adds functionality, the docs should be updated. Put your new functionality into a function with a docstring, and add the feature to the list in *README.rst*.
3. The pull request should work for Python 2.6, 2.7, 3.3, 3.4 and 3.5, and for PyPy. Check [https://travis-ci.org/cdemulde/wwdata/pull\\_requests](https://travis-ci.org/cdemulde/wwdata/pull_requests) and make sure that the tests pass for all supported Python versions.

## 5.4 Tips

To run a subset of tests:

```
$ python -m unittest tests.test_wwdata
```





CREDITS

## 6.1 Development Lead

Chaim De Mulder

## 6.2 Contributors

None yet. All contributions are welcome!

## 6.3 Current funder

Waterboard De Dommel (The Netherlands)



## HISTORY

### 7.1 0.1.0 (2017-09-01)

- First release on PyPI.

The `wwdata` (wastewater data) package is meant to make data analysis, validation and filling of data gaps more streamlined. It contains code to do all this, while also providing simple visualisations of the whole procedure.

The package was (and is) developed in the framework of PhD research, involving the modelling of a full scale wastewater treatment plant (WWTP). Online measurements at the plant are available, but as with all data, is not perfect and therefor needs validation. The gap filling originated from the need to have high-frequency influent data available to run the WWTP model with.



## INDICES AND TABLES

- genindex
- modindex
- search



## PYTHON MODULE INDEX

### W

`wwdata`, [31](#)  
`wwdata.Class_HydroData`, [11](#)  
`wwdata.Class_LabExperimBased`, [20](#)  
`wwdata.Class_LabSensorBased`, [21](#)  
`wwdata.Class_OnlineSensorBased`, [22](#)  
`wwdata.data_reading_functions`, [28](#)  
`wwdata.time_conversion_functions`, [30](#)

# INDEX

## A

absolute\_to\_relative() (in module ww-  
data.Class\_OnlineSensorBased), 27

absolute\_to\_relative() (ww-  
data.Class\_HydroData.HydroData method),  
11

add\_conc() (wwdata.Class\_LabExperimBased.LabExperimBased  
method), 20

add\_to\_filled() (wwdata.Class\_OnlineSensorBased.OnlineSensorBased  
method), 23

add\_to\_meta\_valid() (ww-  
data.Class\_HydroData.HydroData method),  
12

## C

calc\_daily\_average() (ww-  
data.Class\_OnlineSensorBased.OnlineSensorBased  
method), 23

calc\_daily\_profile() (ww-  
data.Class\_HydroData.HydroData method),  
12

calc\_ratio() (wwdata.Class\_HydroData.HydroData  
method), 12

calc\_slope() (wwdata.Class\_LabExperimBased.LabExperimBased  
method), 20

calc\_slopes() (wwdata.Class\_HydroData.HydroData  
method), 13

calc\_total\_proportional() (ww-  
data.Class\_OnlineSensorBased.OnlineSensorBased  
method), 23

check\_filling\_error() (ww-  
data.Class\_OnlineSensorBased.OnlineSensorBased  
method), 23

check\_ph() (wwdata.Class\_LabExperimBased.LabExperimBased  
method), 20

compare\_ratio() (wwdata.Class\_HydroData.HydroData  
method), 13

## D

data\_type (wwdata.Class\_HydroData.HydroData at-  
tribute), 11

data\_type (wwdata.Class\_LabExperimBased.LabExperimBased  
attribute), 20

data\_type (wwdata.Class\_LabSensorBased.LabSensorBased  
attribute), 21

data\_type (wwdata.Class\_OnlineSensorBased.OnlineSensorBased  
attribute), 22

drop\_index\_duplicates() (ww-  
data.Class\_HydroData.HydroData method),  
13

drop\_index\_duplicates() (ww-  
data.Class\_OnlineSensorBased.OnlineSensorBased  
method), 24

drop\_peaks() (in module ww-  
data.Class\_OnlineSensorBased), 27

drop\_peaks() (wwdata.Class\_LabSensorBased.LabSensorBased  
method), 22

## E

experiment\_tag (wwdata.Class\_HydroData.HydroData  
attribute), 11

experiment\_tag (wwdata.Class\_LabExperimBased.LabExperimBased  
attribute), 20

experiment\_tag (wwdata.Class\_LabSensorBased.LabSensorBased  
attribute), 21

experiment\_tag (wwdata.Class\_OnlineSensorBased.OnlineSensorBased  
attribute), 22

## F

fill\_index() (wwdata.Class\_HydroData.HydroData  
method), 13

fill\_missing\_correlation() (ww-  
data.Class\_OnlineSensorBased.OnlineSensorBased  
method), 24

fill\_missing\_daybefore() (ww-  
data.Class\_OnlineSensorBased.OnlineSensorBased  
method), 25

fill\_missing\_interpolation() (ww-  
data.Class\_OnlineSensorBased.OnlineSensorBased  
method), 25

fill\_missing\_model() (ww-  
data.Class\_OnlineSensorBased.OnlineSensorBased  
method), 26



- fill\_missing\_ratio() (wwdata.Class\_OnlineSensorBased.OnlineSensorBased method), 26
- fill\_missing\_standard() (wwdata.Class\_OnlineSensorBased.OnlineSensorBased method), 27
- find\_and\_replace() (in module wwdata.data\_reading\_functions), 28
- find\_nearest\_time() (in module wwdata.Class\_OnlineSensorBased), 28
- ## G
- get\_absolute\_time() (in module wwdata.time\_conversion\_functions), 30
- get\_avg() (wwdata.Class\_HydroData.HydroData method), 13
- get\_correlation() (wwdata.Class\_HydroData.HydroData method), 14
- get\_highs() (wwdata.Class\_HydroData.HydroData method), 14
- get\_std() (wwdata.Class\_HydroData.HydroData method), 14
- go.WEST() (in module wwdata.Class\_OnlineSensorBased), 28
- ## H
- head() (wwdata.Class\_HydroData.HydroData method), 14
- hours() (wwdata.Class\_LabExperimBased.LabExperimBased method), 21
- HydroData (class in wwdata.Class\_HydroData), 11
- ## I
- in\_out() (wwdata.Class\_LabExperimBased.LabExperimBased method), 21
- index() (wwdata.Class\_HydroData.HydroData method), 15
- ## J
- join\_files() (in module wwdata.data\_reading\_functions), 29
- ## L
- LabExperimBased (class in wwdata.Class\_LabExperimBased), 20
- LabSensorBased (class in wwdata.Class\_LabSensorBased), 21
- list\_files() (in module wwdata.data\_reading\_functions), 29
- ## M
- make\_datetime() (in module wwdata.time\_conversion\_functions), 30
- make\_month\_day\_array() (in module wwdata.time\_conversion\_functions), 31
- moving\_average\_filter() (wwdata.Class\_HydroData.HydroData method), 15
- moving\_slope\_filter() (wwdata.Class\_HydroData.HydroData method), 15
- ## O
- OnlineSensorBased (class in wwdata.Class\_OnlineSensorBased), 22
- ## P
- plot() (wwdata.Class\_LabExperimBased.LabExperimBased method), 21
- plot\_analysed() (wwdata.Class\_HydroData.HydroData method), 16
- ## R
- read\_mat() (in module wwdata.data\_reading\_functions), 29
- removal() (wwdata.Class\_LabExperimBased.LabExperimBased method), 21
- remove\_empty\_lines() (in module wwdata.data\_reading\_functions), 29
- replace() (wwdata.Class\_HydroData.HydroData method), 16
- ## S
- savgol() (wwdata.Class\_HydroData.HydroData method), 16
- set\_index() (wwdata.Class\_HydroData.HydroData method), 17
- set\_tag() (wwdata.Class\_HydroData.HydroData method), 17
- set\_time\_unit() (wwdata.Class\_HydroData.HydroData method), 17
- set\_units() (wwdata.Class\_HydroData.HydroData method), 17
- simple\_moving\_average() (wwdata.Class\_HydroData.HydroData method), 17
- sort\_data() (in module wwdata.data\_reading\_functions), 29
- ## T
- tag\_doubles() (wwdata.Class\_HydroData.HydroData method), 18
- tag\_extremes() (wwdata.Class\_HydroData.HydroData method), 18
- tag\_nan() (wwdata.Class\_HydroData.HydroData method), 19
- tail() (wwdata.Class\_HydroData.HydroData method), 19

time\_unit (wwdata.Class\_HydroData.HydroData attribute), 11

time\_unit (wwdata.Class\_LabExperimBased.LabExperimBased attribute), 20

time\_unit (wwdata.Class\_LabSensorBased.LabSensorBased attribute), 21

time\_unit (wwdata.Class\_OnlineSensorBased.OnlineSensorBased attribute), 22

timedata\_column (wwdata.Class\_HydroData.HydroData attribute), 11

timedata\_column (wwdata.Class\_LabExperimBased.LabExperimBased attribute), 20

timedata\_column (wwdata.Class\_LabSensorBased.LabSensorBased attribute), 21

timedata\_column (wwdata.Class\_OnlineSensorBased.OnlineSensorBased attribute), 22

timedelta\_to\_abs() (in module wwdata.time\_conversion\_functions), 31

to\_datetime() (wwdata.Class\_HydroData.HydroData method), 19

to\_datetime\_singlevalue() (in module wwdata.time\_conversion\_functions), 31

to\_days() (in module wwdata.time\_conversion\_functions), 31

to\_float() (wwdata.Class\_HydroData.HydroData method), 19

to\_hours() (in module wwdata.time\_conversion\_functions), 31

to\_minutes() (in module wwdata.time\_conversion\_functions), 31

to\_seconds() (in module wwdata.time\_conversion\_functions), 31

total\_seconds() (in module wwdata.Class\_HydroData), 19

total\_seconds() (in module wwdata.Class\_OnlineSensorBased), 28

## U

units (wwdata.Class\_HydroData.HydroData attribute), 11

units (wwdata.Class\_LabExperimBased.LabExperimBased attribute), 20

units (wwdata.Class\_LabSensorBased.LabSensorBased attribute), 22

units (wwdata.Class\_OnlineSensorBased.OnlineSensorBased attribute), 23

## V

vlookup\_day() (in module wwdata.Class\_OnlineSensorBased), 28

## W

write() (wwdata.Class\_HydroData.HydroData method), 19

write\_to\_WEST() (in module wwdata.data\_reading\_functions), 30

wwdata (module), 31

wwdata.Class\_HydroData (module), 11

wwdata.Class\_LabExperimBased (module), 20

wwdata.Class\_LabSensorBased (module), 21

wwdata.Class\_OnlineSensorBased (module), 22

wwdata.data\_reading\_functions (module), 28

wwdata.time\_conversion\_functions (module), 30