# Statistical functions

## Statistical functions

## Box plot

The box plot function will calculate a set of statistical values for each of the selected timeseries and displays these values in the chart as a Box and Whisker Plot.

**Time Series Display**

Functions: Box plot

**Box plot**
**28-02-2010 11:50:00 - 28-02-2010 14:15:00**

| Statistics | B Neerslaghoeve (mm) 0306-AME-01 r0000518 | C Neersla (mm) 0306-A r000051 |
|---|---|---|
| Minimum outlier | 0.004 | |
| Minimum regular value | 0.004 | |
| 25% | 0.021 | |
| Median | 0.044 | |
| Mean | 0.052 | |
| 75% | 0.075 | |
| Maximum regular value | 0.131 | |
| Maximum outlier | 0.131 | |

Neerslaghoeveelheid 5 min (mm)

- Maximum regular value
- 75% percentile
- Mean
- Median
- 25% percentile
- Minimum regular value

Neerslaghoeveelheid 5 m...    Neerslaghoeveelheid 5 m...    Neerslaghoeveelheid 5 m...

■ Neerslaghoeveelheid 5 min 0306-AME-01  ■ Neerslaghoeveelheid 5 min 0306-AME-02
■ Neerslaghoeveelheid 5 min 0306-AME-03

Close    Help

Current system time: 09-03-2010 08:00 GMT    08:26:37 GMT    10:26:37 CEST    Display time: 09-03-2010 08:00 GMT    Stand alone    Last refresh: 10:25:50 CEST    N...

In the table the following list of statical values is shown for each of the timeseries:

**Statistical variables**
*Minimum outlier*: minimum outlier value for selection.
*Minimum regular value*: minimum value for selection that is not defined as an outlier. Also known as the Whiskers
*25%*: 25th percentile
*Median*: 50th percentile
*Mean*: average value for selection
*75%*: 75th percentile*Maximum regular value*: maximum value for selection that is not defined as an outlier. Also known as the Whiskers*Maximum regular value*: maximum value for selection that is not defined as an outlier. Also known as the Whiskers
*Maximum outlier*: maximum outlier value for selection.

The chart shows the same values as the table, hover the chart also can include some extra values that are not shown in the table. These are the outliers and the far-out indicators.
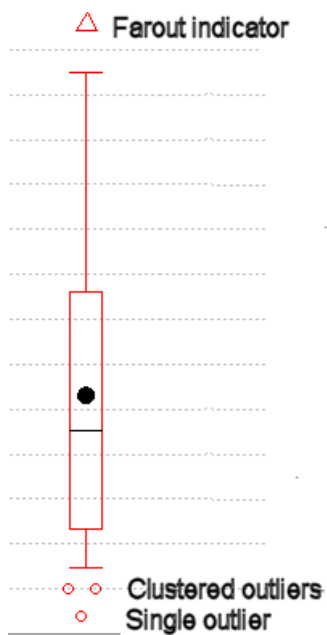
*Far-out indicator*: indicates that there are values that lie outside the plotted range of the axis.
*Single outlier:* a single outlier value.
*Clustered outlier*: multiple outliers that are located too close together to be plotted separately.

**Outliers**: cases where the values are between 1.5 and 3 box-lengths from the 75th percentile or 25th percentile.
**Farout values**: cases where the values are more than 3 box-lengths from the 75th percentile or 25th percentile.

Config example:

```
<statisticalFunctions>
              <statisticalFunction function="boxPlot"/>
</statisticalFunctions>
```

# Calendar aggregation

Aggregation by calendar day (00.00h-24.00h)
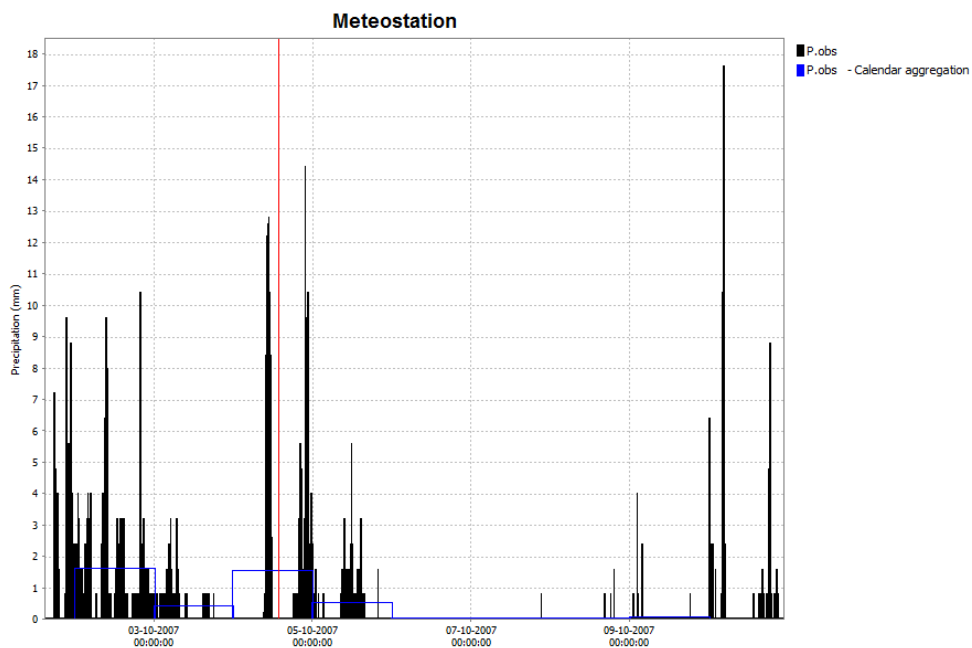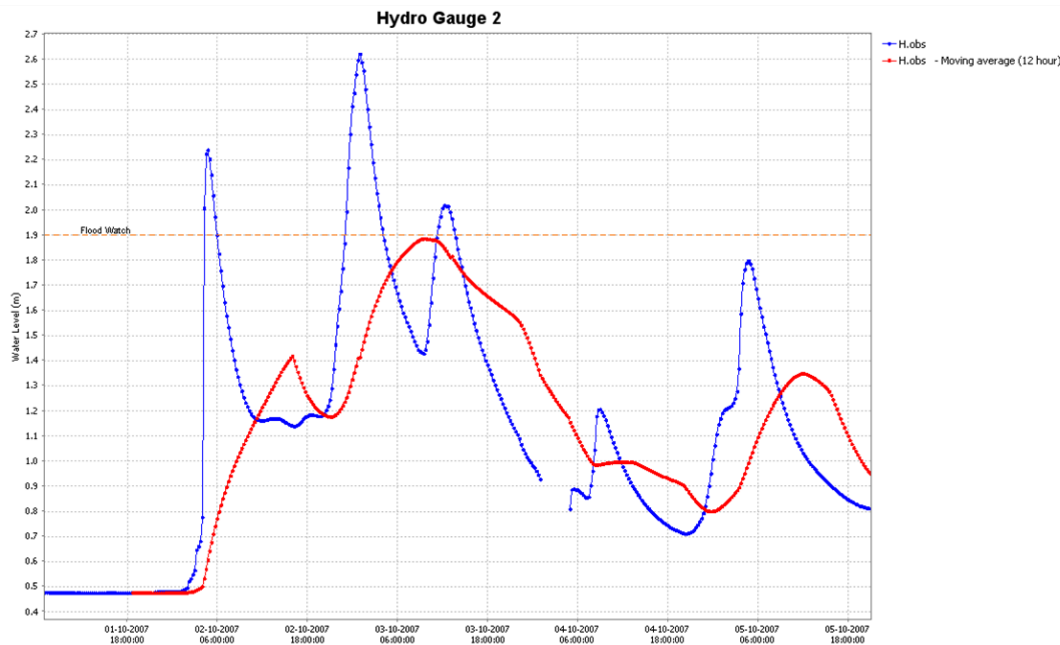


Config example:

```
<statisticalFunctions>
<statisticalFunction function="calendarAggregation">
                    <movingAccumulationTimeSpan unit="week" multiplier="1"/>
                    <movingAccumulationTimeSpan unit="week" multiplier="2"/>
                    <movingAccumulationTimeSpan unit="week" multiplier="4"/>
                    <movingAccumulationTimeSpan unit="week" multiplier="12"/>
                    <movingAccumulationTimeSpan unit="day" multiplier="365"/>
            </statisticalFunction>
```

In addition to the calender aggregation, there is also the option for accumulationAggregation and relativeAggregation. There is also a very similar function called accumulationInterval where the accumulation is displayed at every timestep (for configuration see example above).

# Moving average

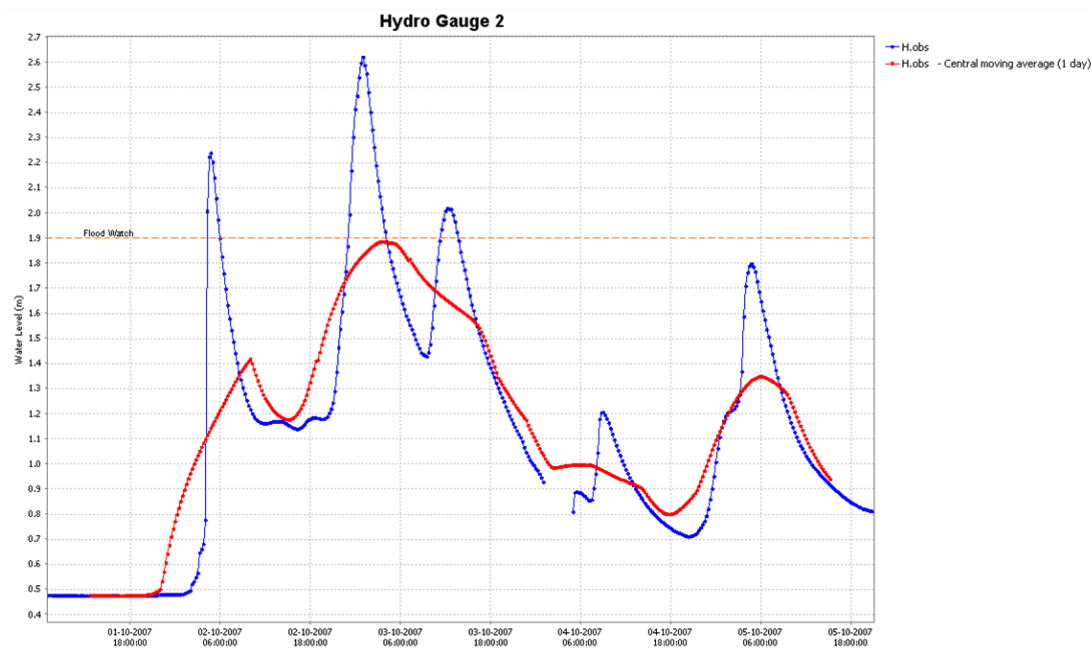Moving average where the value is stamped at the end of each averaging period.



Config example:

```
<statisticalFunctions>
            <statisticalFunction function="movingAverage" ignoreMissings="true">
                    <movingAccumulationTimeSpan unit="hour" multiplier="1"/>
                    <movingAccumulationTimeSpan unit="hour" multiplier="3"/>
                    <movingAccumulationTimeSpan unit="hour" multiplier="6"/>
                    <movingAccumulationTimeSpan unit="hour" multiplier="12"/>
            </statisticalFunction>
```

# Central moving average

Moving average where the value is stamped in the middle of each averaging period.

Hydro Gauge 2

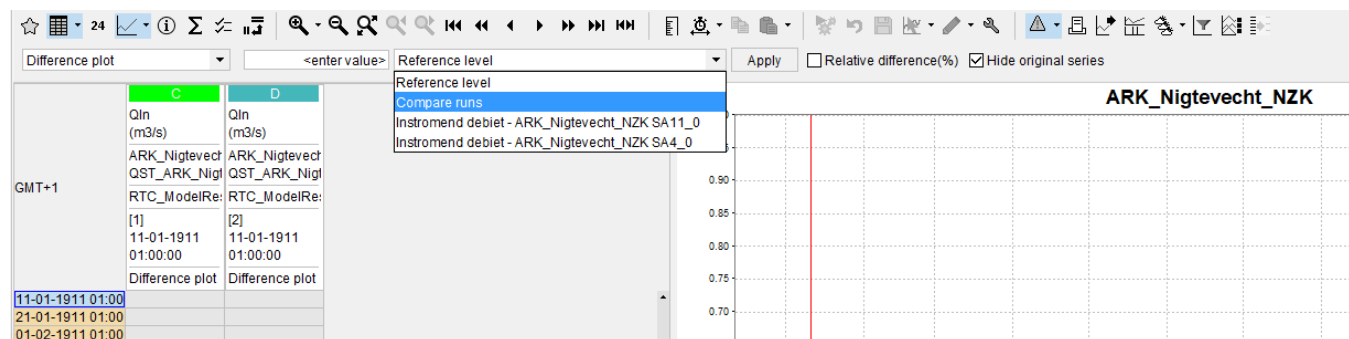Config example:

```
<statisticalFunctions>
            <statisticalFunction function="centralMovingAverage" ignoreMissings="true">
                    <movingAccumulationTimeSpan unit="hour" multiplier="1"/>
                    <movingAccumulationTimeSpan unit="hour" multiplier="3"/>
                    <movingAccumulationTimeSpan unit="hour" multiplier="6"/>
                    <movingAccumulationTimeSpan unit="hour" multiplier="12"/>
            </statisticalFunction>
```

# Differences

Displays the difference between two selected series

Config example:

```
<statisticalFunction function="differences"/>
```



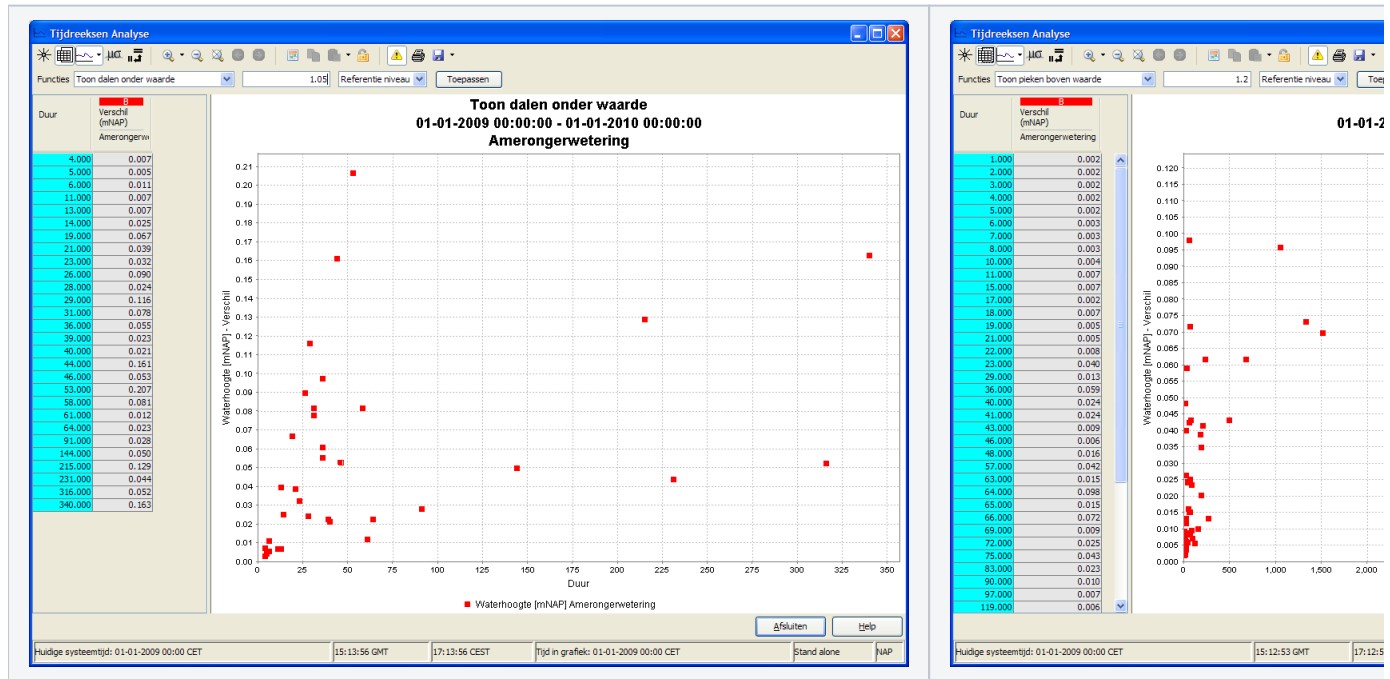If you configure Differences statistical function, the following options will become available in plot:

- Reference level.
  This option is always available. It compares the values of the TimeSeriesSet to a value you can specify by typing into the box marked <enter value>, and displays the difference: TimeSeriesSet value - reference value.
- Compare runs (available since 2020.01)
  This option is only available if you have opened/selected at least 2 different runs for each timeSeriesSet in the display. The number can be greater than 2, and does not need to be the same number for all timeSeries. (For example 2 different runs for parameter A, and 3 for B). As reference level, the oldest of this runs is selected, and the difference compared to this run is displayed. (Run value - reference run value)

- Select timeSeries as reference level
  This option is available if you have at least two different time series open in the plot. They can be different runs, but they do not need to be. The difference between the Time series and reference time series will be displayed. (TimeSeries value - reference TimeSeries value)

If the checkbox "Relative difference(%)" is checked, the relative difference will be displayed (in %). This is available for all three options.  100*(TimeSeries value - reference TimeSeries value)/reference TimeSeries value.  If the reference value is 0, N.a.N will be the result.

## Display lows below value & Display peaks above value



Shows all peak heights or dip depths and duration above or below a reference level. This level is set by default but can be altered in the toolbar. After adjusting, press the apply button to recalculate the peaks or dips. . From now on we will only mention the peaks and how values are above a level, but this functionality is symmetrical for the valleys and values being below a level. The default value for the reference level is set by determining the 'low' areas according to the maximum available value of the input time series array.

After setting the reference level, every continues series of values above that level is considered a peak and will be depicted as a dot in the XY plot. The X axis shows what the duration of this peak is, with this value corresponding to the amount of time entries that make up that peak. The Y axis shows the average height of the peak above the reference level. The parameter used on the Y axis is equal to that of the selected time series.

When hovering above a point in the XY plot, a tooltip will appear. Here is an overview of the tooptip format with some of the elements being optional:

- |Parameter Id| |Location Id| (|Peak duration|, |Average peak height|), Max Difference = |Max peak height|, Percentile Difference = |Percentile peak height| at time |Peak time|
- Parameter Id: Paramter Id from selected timeseries
- Location Id: Location Id from selected timeserires
- Peak duration: Number of time entries in peak
- Average peak height: Average of all differences between reference level and the time entries value
- Max peak height (optional): Maximum of all differences between reference level and the time entries value
- Percentile peak height (optional): Relation between the average and max peak height
- Peak time: Time of the first time entry in the peak

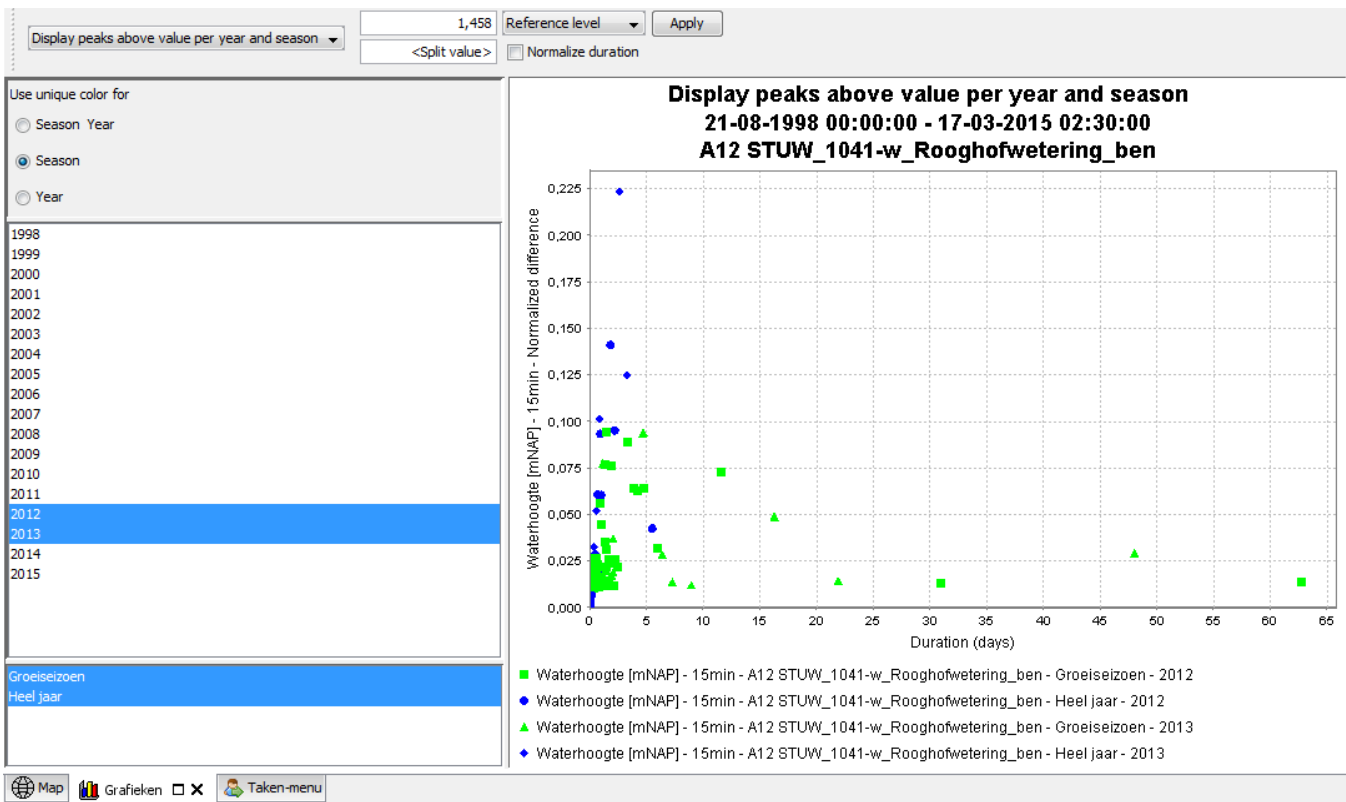Config example:

```
<statisticalFunction function="showLowsBelow"/>
<statisticalFunction function="showPeaksAbove"/>
```
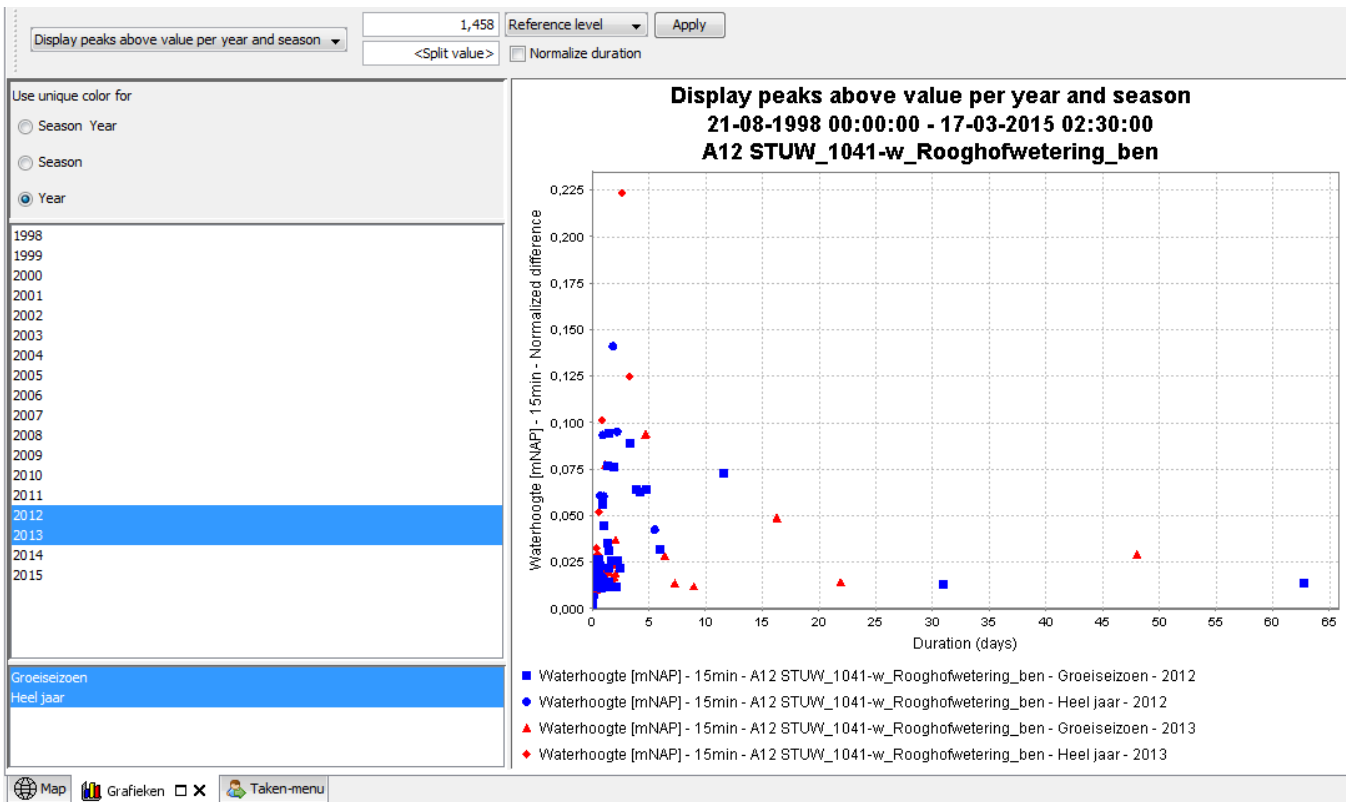
## Display lows below value & Display peaks above value (per year and season)

(Since 2016.02) The same as the function above but then with the extension of selecting multiple different years and seasons, plus a choice of how the different seasons and years should be uniquely colored.
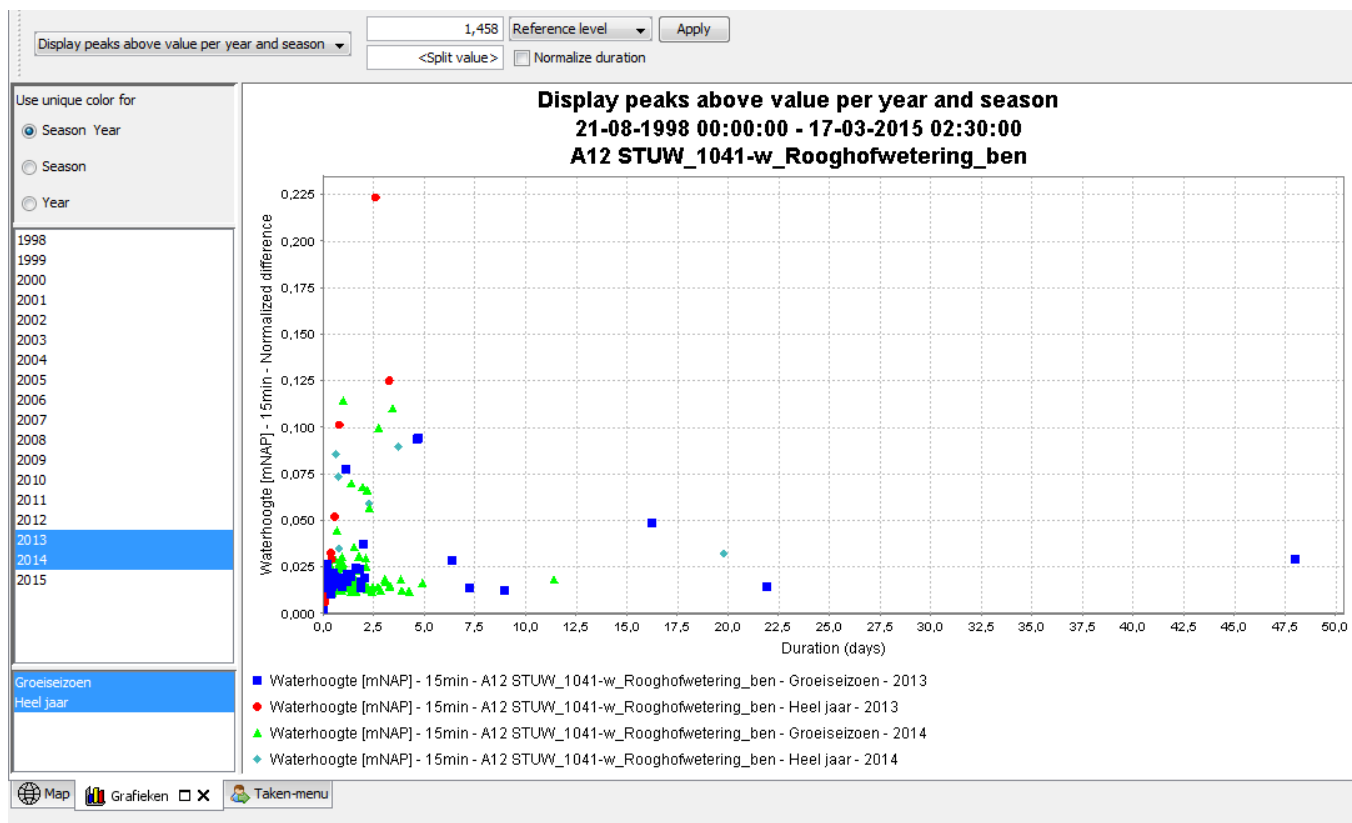
Unique color per season

Unique color per year:



Unique color per year and season:

**Example config**
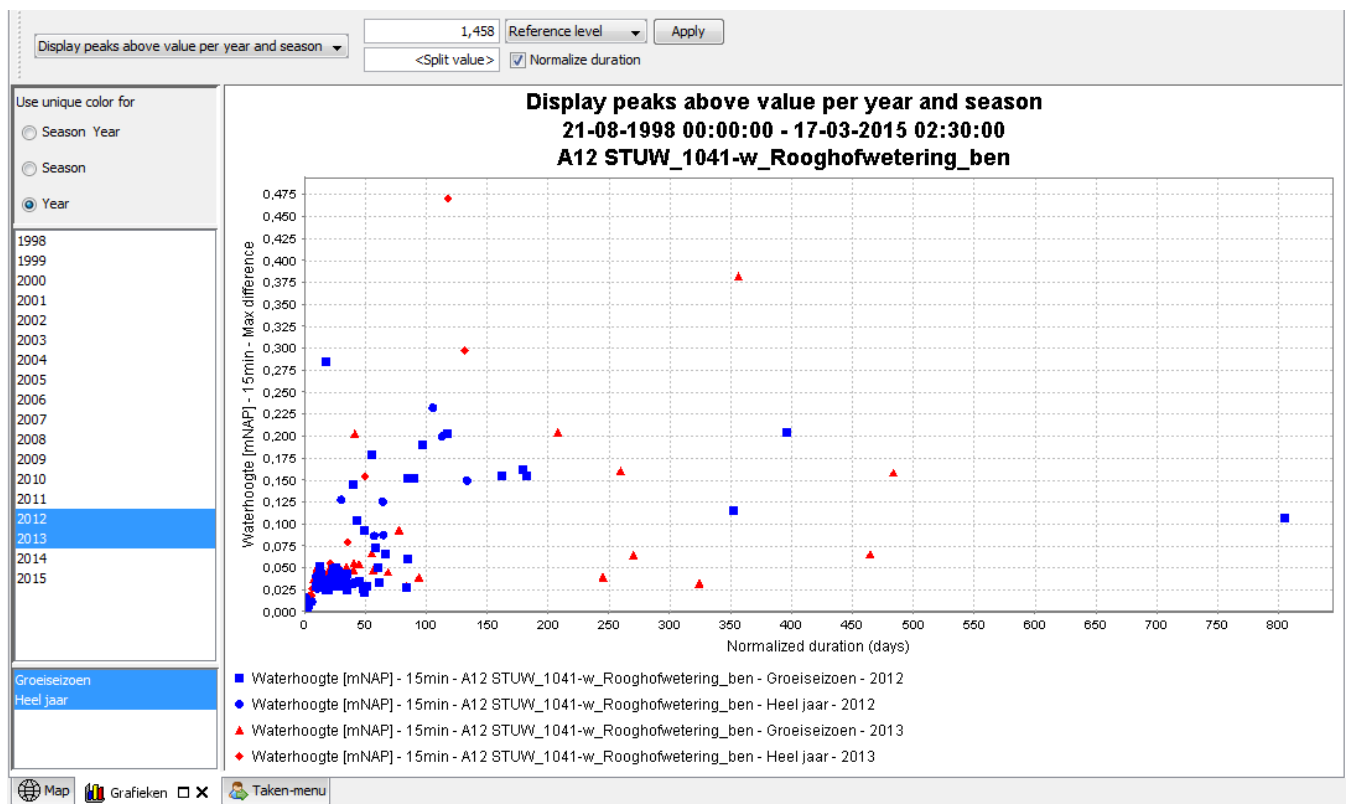
```
<statisticalFunction function="historicalShowPeaksAbove">
        <season startMonthDay="--01-01" endMonthDay="--12-31" label="Year" color="black"/>
        <season startMonthDay="--04-01" endMonthDay="--09-30" label="Grow Season" color="gray"/>
</statisticalFunction>
<statisticalFunction function="historicalShowLowsBelow">
        <season startMonthDay="--01-01" endMonthDay="--12-31" label="Year" color="black"/>
        <season startMonthDay="--04-01" endMonthDay="--09-30" label="Grow Season" color="gray"/>
</statisticalFunction>
```
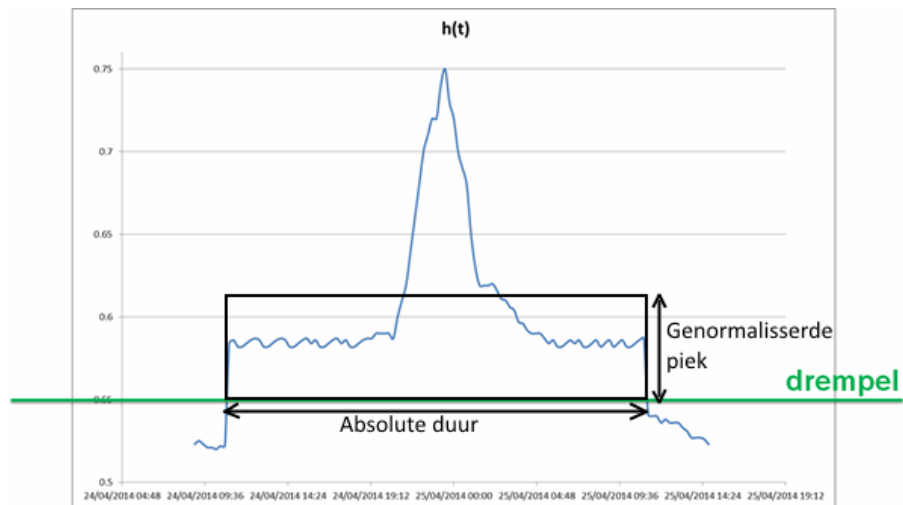
## Normalize duration

(Since 2016.02)  for the peaks above/below functions a checkbox is added to normalize over duration. This changes The y-axis from Normalized difference to Max difference and the X-axis from (total) duration the Normalized duration.
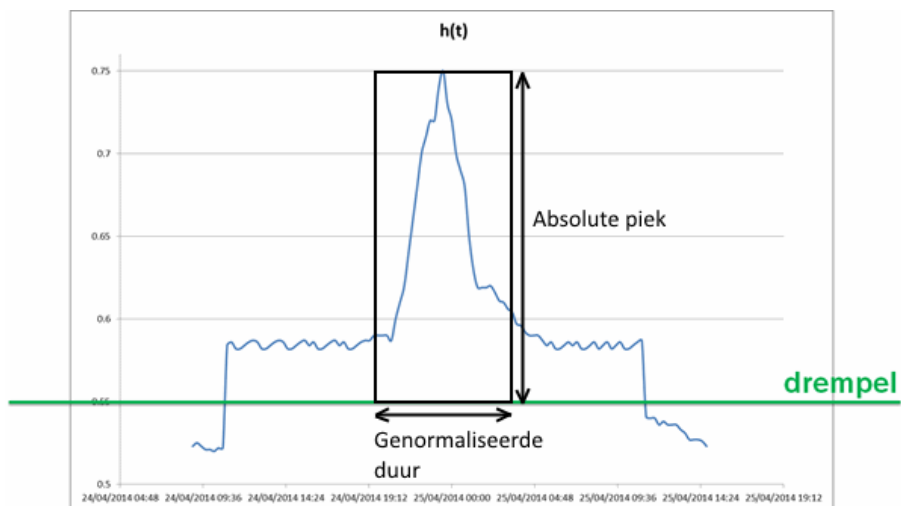
The differences between normalization can best be shown in a pictures
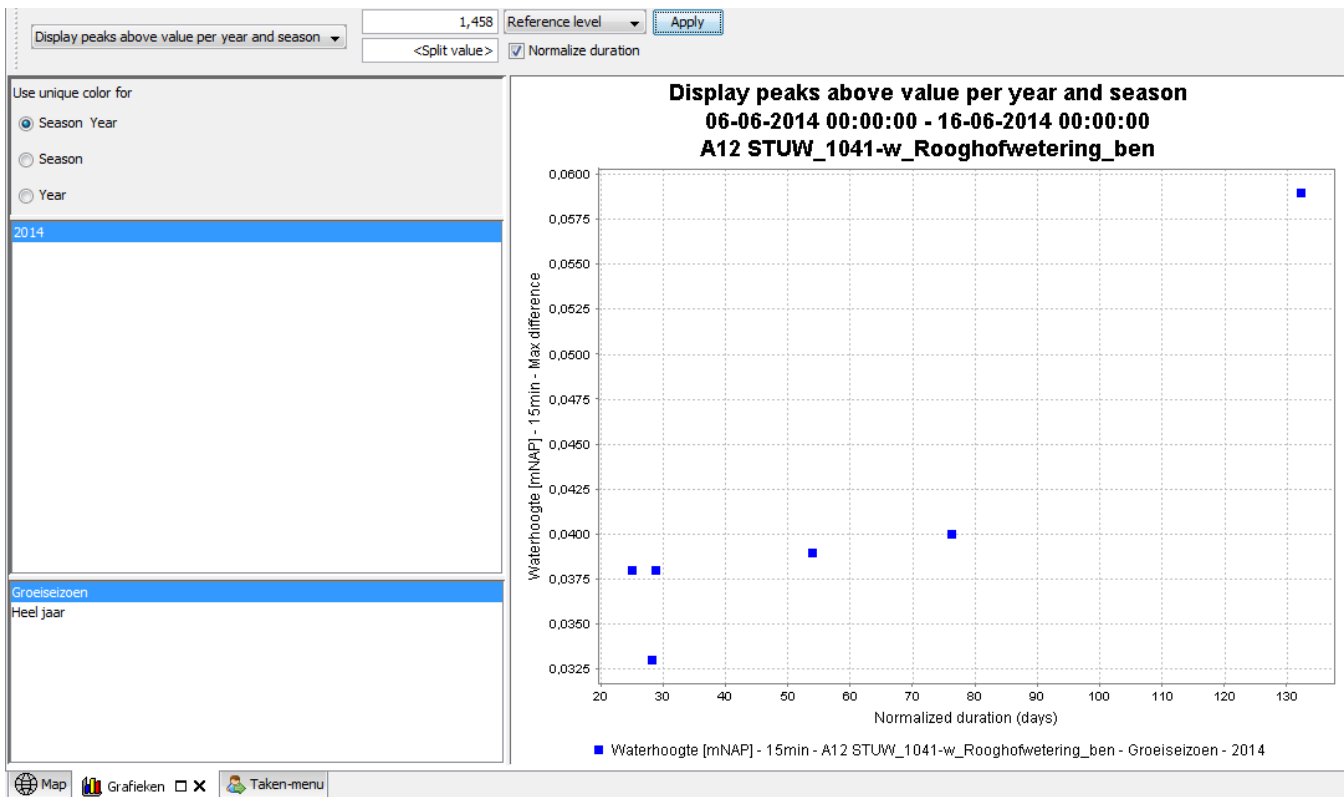
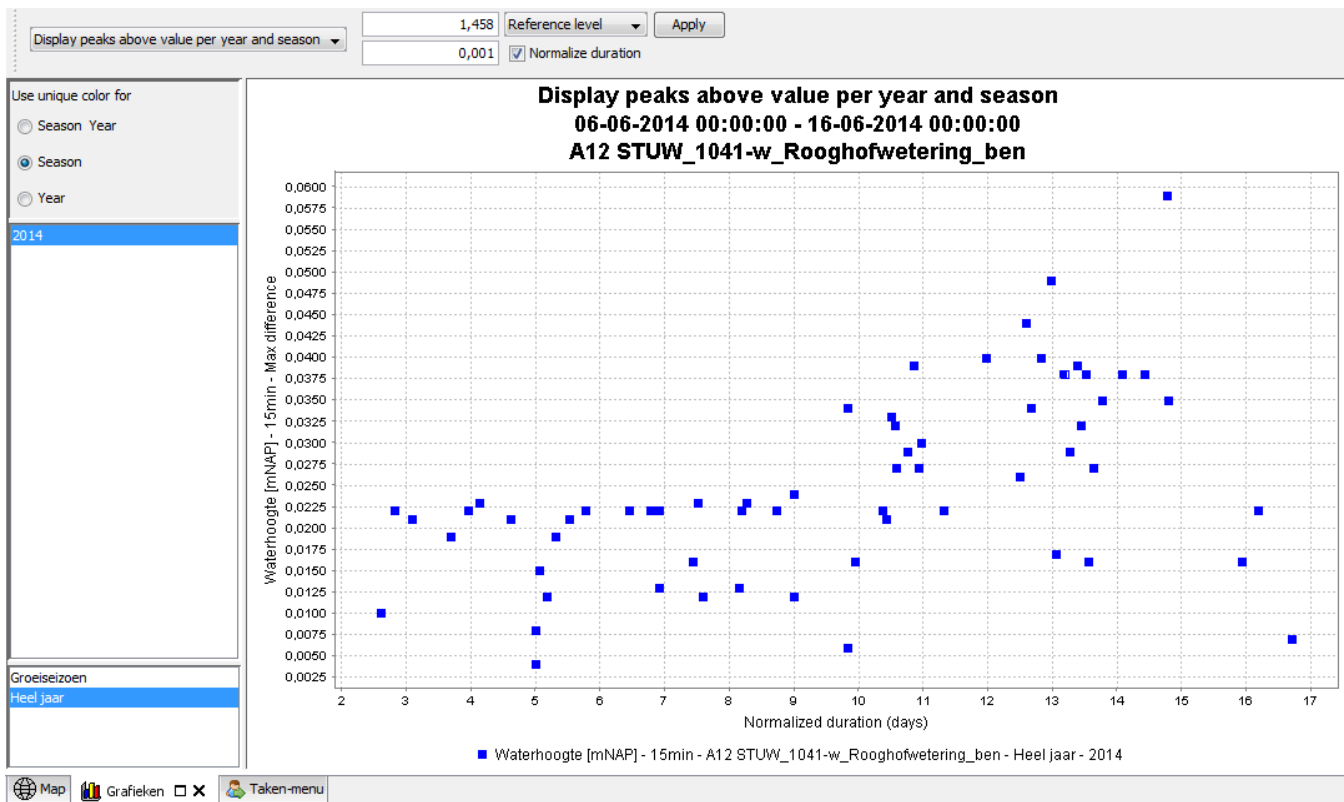Normalized difference:



Normalized duration:

## Split peaks

(Since 2016.02) for the peaks above/below functions a possibility is added to split the peaks when they have multiple local maxima.
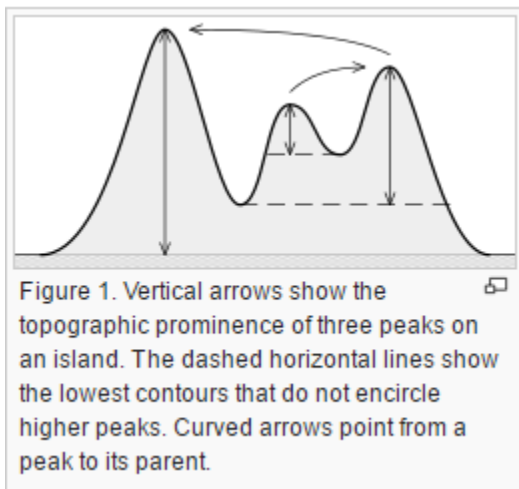
No split:



Split when local maxima have a low between them with a value difference of at least 0,001. As can be seen this results in a lot more separate peaks.

For the value difference the measure Prominence is used, which is an official measure for determining separate peaks.
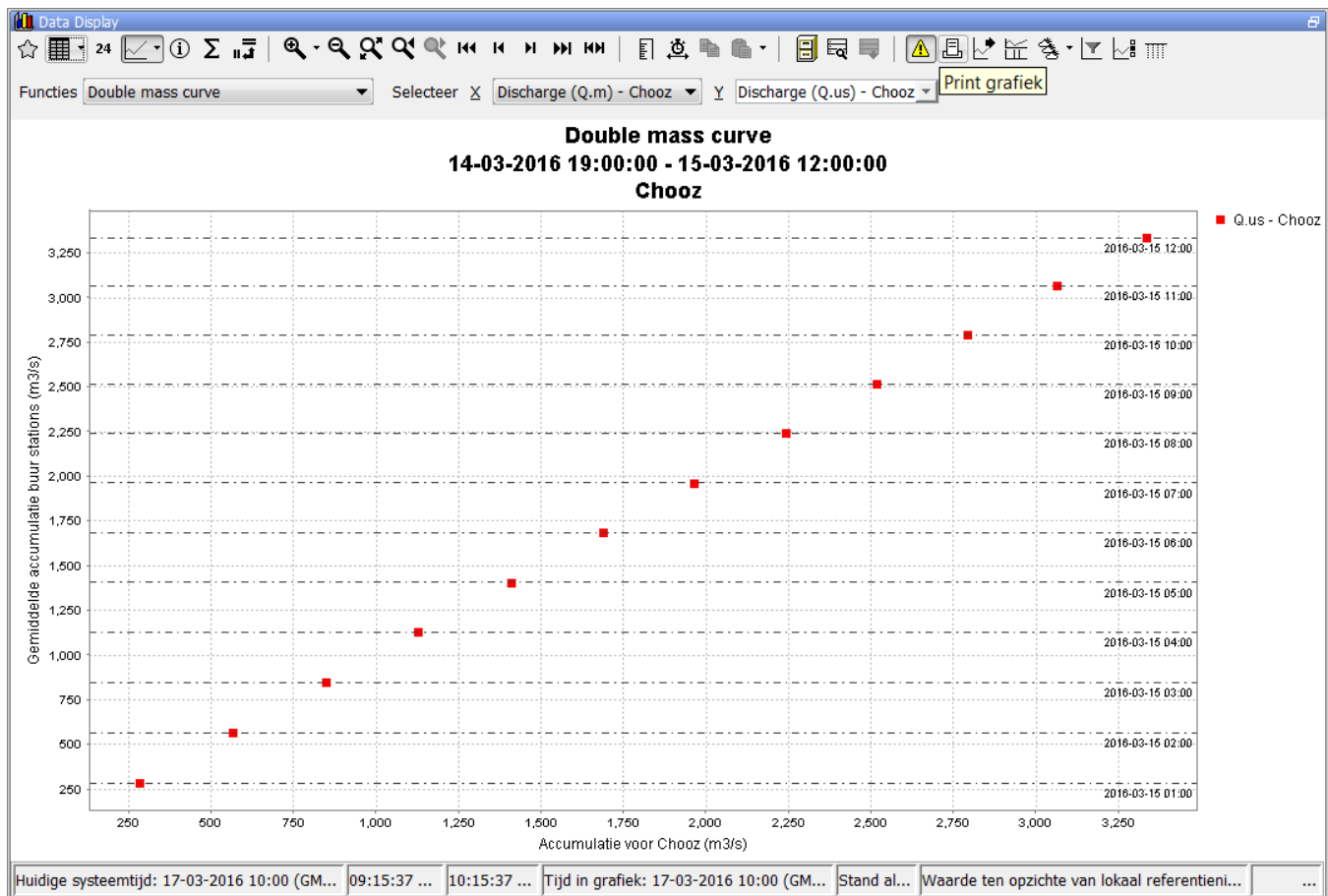


Figure 1. Vertical arrows show the topographic prominence of three peaks on an island. The dashed horizontal lines show the lowest contours that do not encircle higher peaks. Curved arrows point from a peak to its parent.

Source: https://en.wikipedia.org/wiki/Topographic_prominence

# Double Mass curve

Isused in the study of the consistency and long-term trend test of hydrometeo--rological data. This method was first used to analyze the consistency of precipitation data in Susquehanna watershed United States by Merriam at 1937 (Merriam, 1937), and Searcy made a theoretical explanation

of it (Searcy, et al., 1960). The theory of the double-mass curve is based on the fact that a plot of the two cumulative quantities during the same period exhibits a straight line so long as the proportionality between the two remains unchanged, and the slope of the line represents the proportionality. This method can smooth a time series and suppress random elements in the series, and thus show the main trends of the time series. In recent 30 years, Chinese scholars analyzed the effect of soil and water conservation measures and land use/ cover changes on runoff and sediment using double mass curve method, and have achieved good results (Mu, et al., 2010). In this study, double-mass curves of precipitation vs streamflow and precipitation vs sediment are plotted for the two contrastive periods to estimate changes in regression slope (proportionality) to quantify the overall efficiency of soil conservation measures before and after transition years.
An example of the doubleMassCurve plot is given here:

Config example:
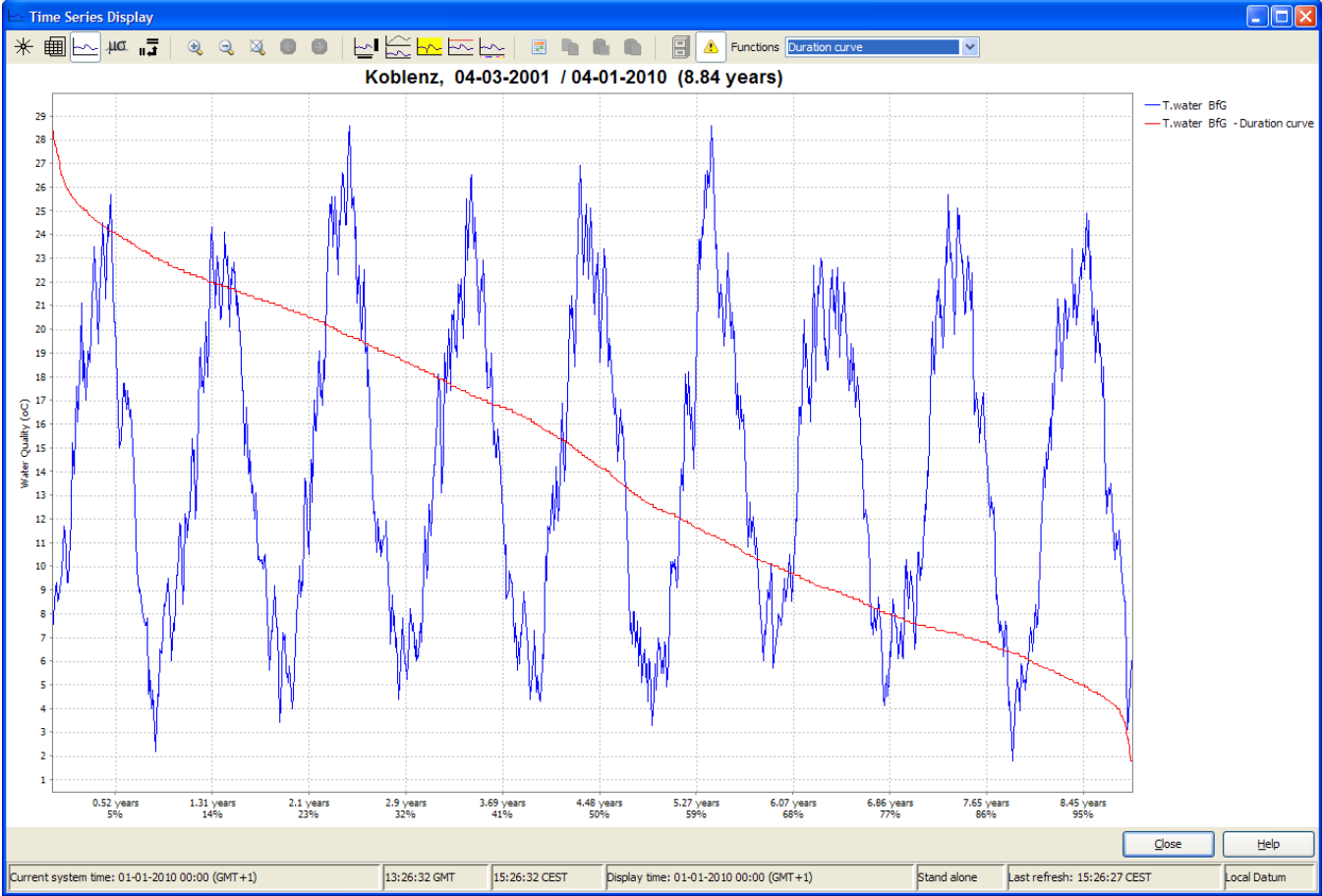
```
<statisticalFunction function="doubleMassCurve"/>
```
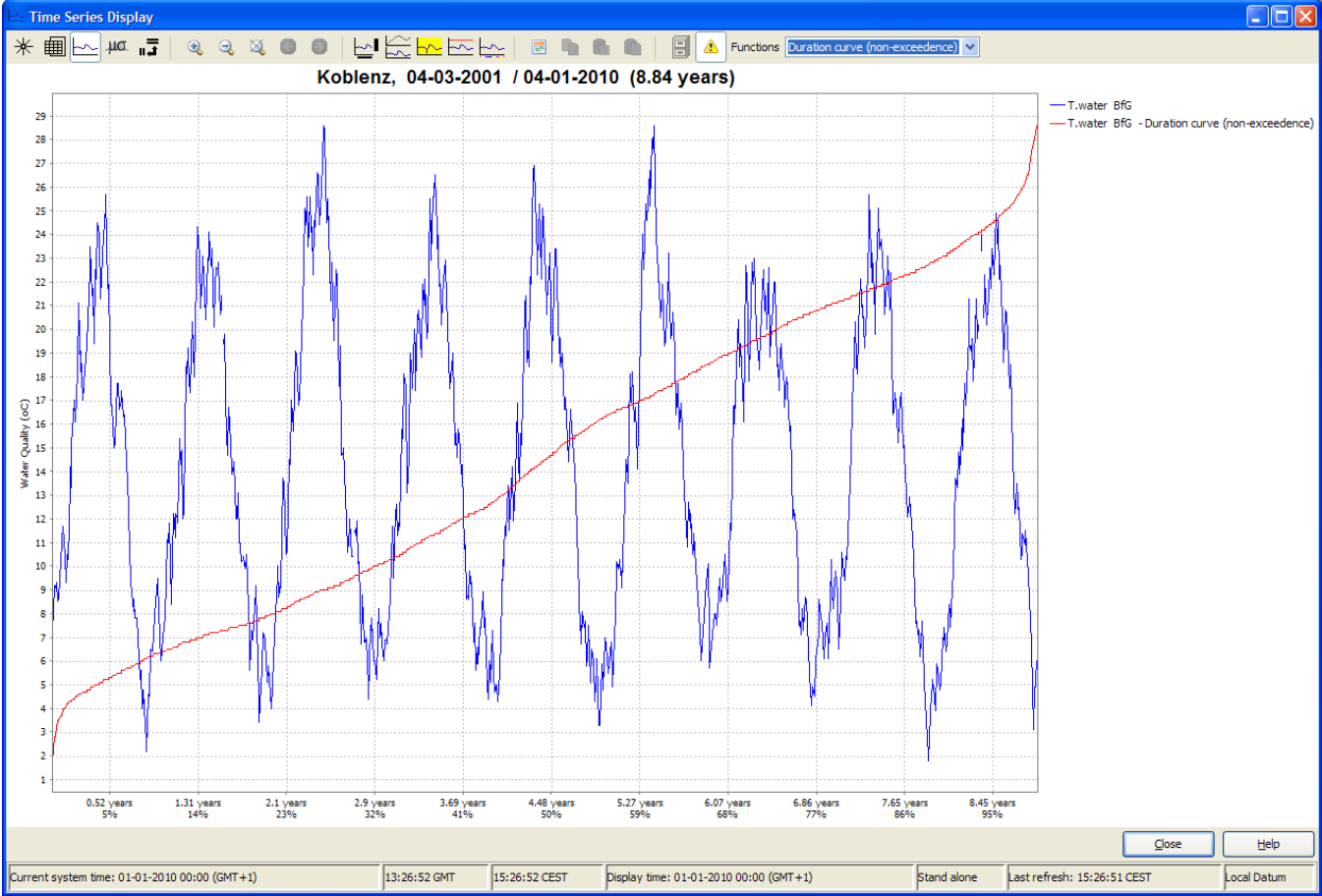
## Duration curve

A convenient way to show the variation of hydrological quantities through time may be done by means of duration curves. For the selected time period the values of the selected quantity are sorted descending (durationExceedence) or ascending (durationNonExceedence). When the duration curve is plotted in the timeseries display, the x-axis will show the entire length in time of the selected view period. Percentages are shown as duration with respect to the entire chosen view period.

In the configuration of this statistical function there is the option to ommit missing values which may occur in the selected view period. If this option is set to true, all entries with missing values will be disregarded before the duration curve is calculated. If this option is not defined (default) or is set to false, missing values will be added to the the end of the array. In this case the plotted duration curve will never reach the 100%.
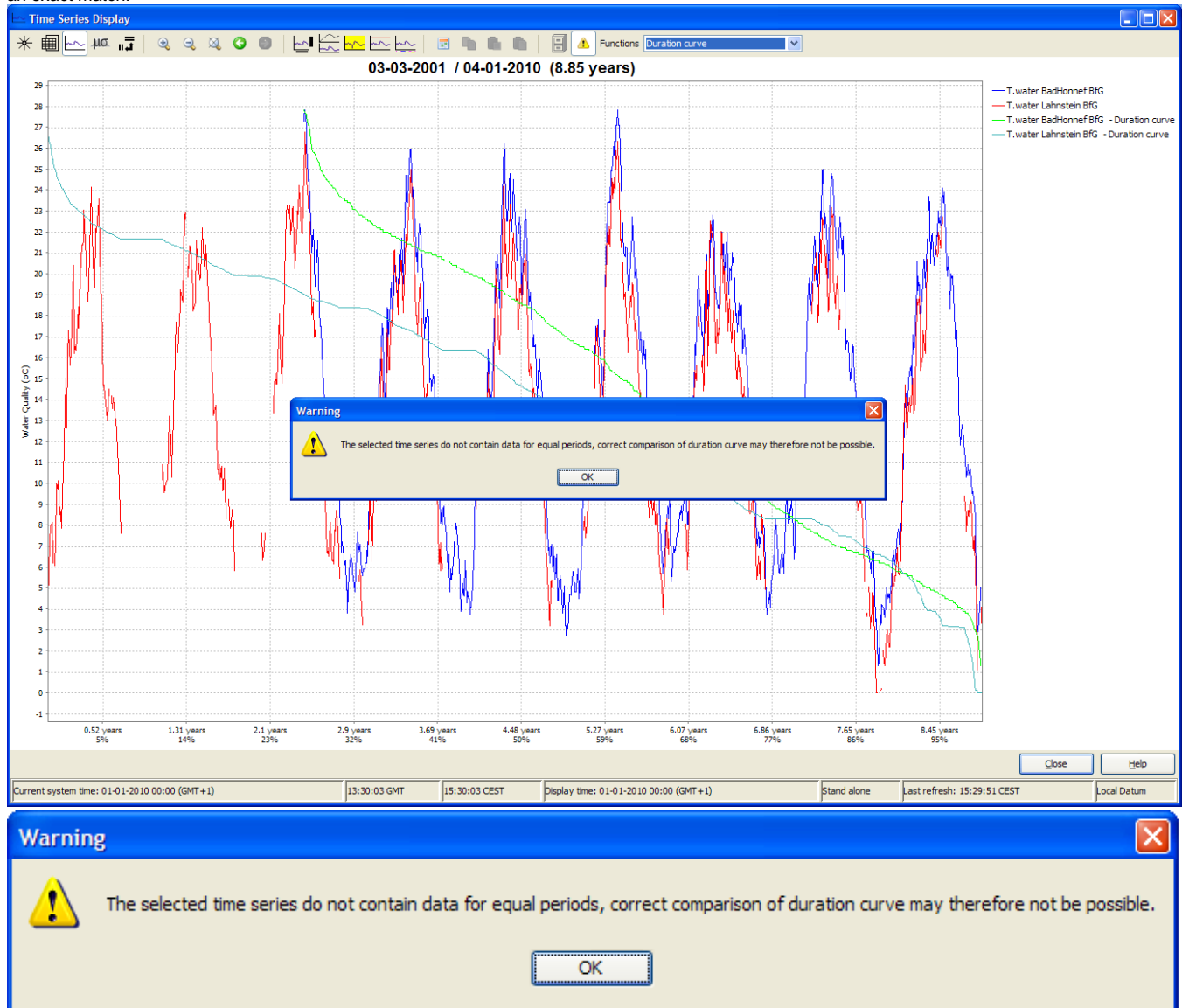
An example of the durationExceedence plot is given here:

An example of the durationNonExceedence plot is given here:

When selecting more than one location it could occur that the view periods of these selected timeseries do not cover the same period in time. In this case it is difficult to make a correct comparison of the calculated duration curves because they are analysed on different periods in time. A warning message will be given in order to ensure that the user is aware of this. The pop-up message will be shown each time the user zooms in or out until all view periods are an exact match.





Config example:

```
<statisticalFunction function="durationExceedence" ignoreMissings="true"/>
<statisticalFunction function="durationNonExceedence" ignoreMissings="true"/>
```

# Elevation

The Elevation statistical function show parameter values against the elevation of the location on one particular time. The time for which the values are shown can be altered with a slider bar.
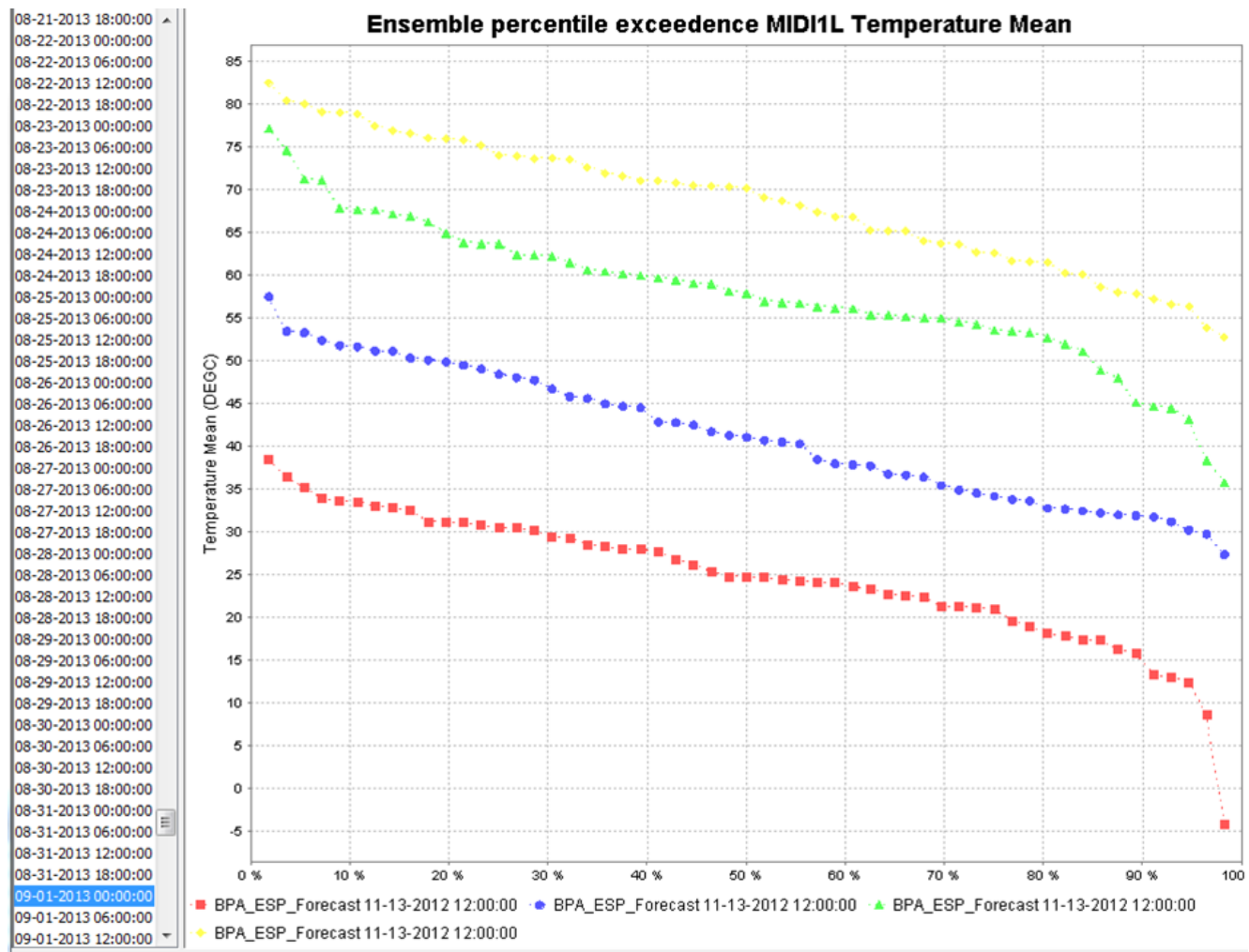
An example of the elevation plot is given here:

Configuration example:

```
<statisticalFunction function="elevation"/>
```

# Ensemble Percentile Exceedence

Plots for the selected time stamp(s) each member of an ensemble along the horzontal axis, sorted by value. The result is exceedence diagrams for an ensemble forecast at a selected timestamp(s) within the timeseries.
The function allows multiple locations and/or multiple time stamps and/or multiple forecasts in one diagram. If multiple forecats include the same selected timestamp, the legend only can distinguish between these forecasts if a taskdescription is used during job submission.
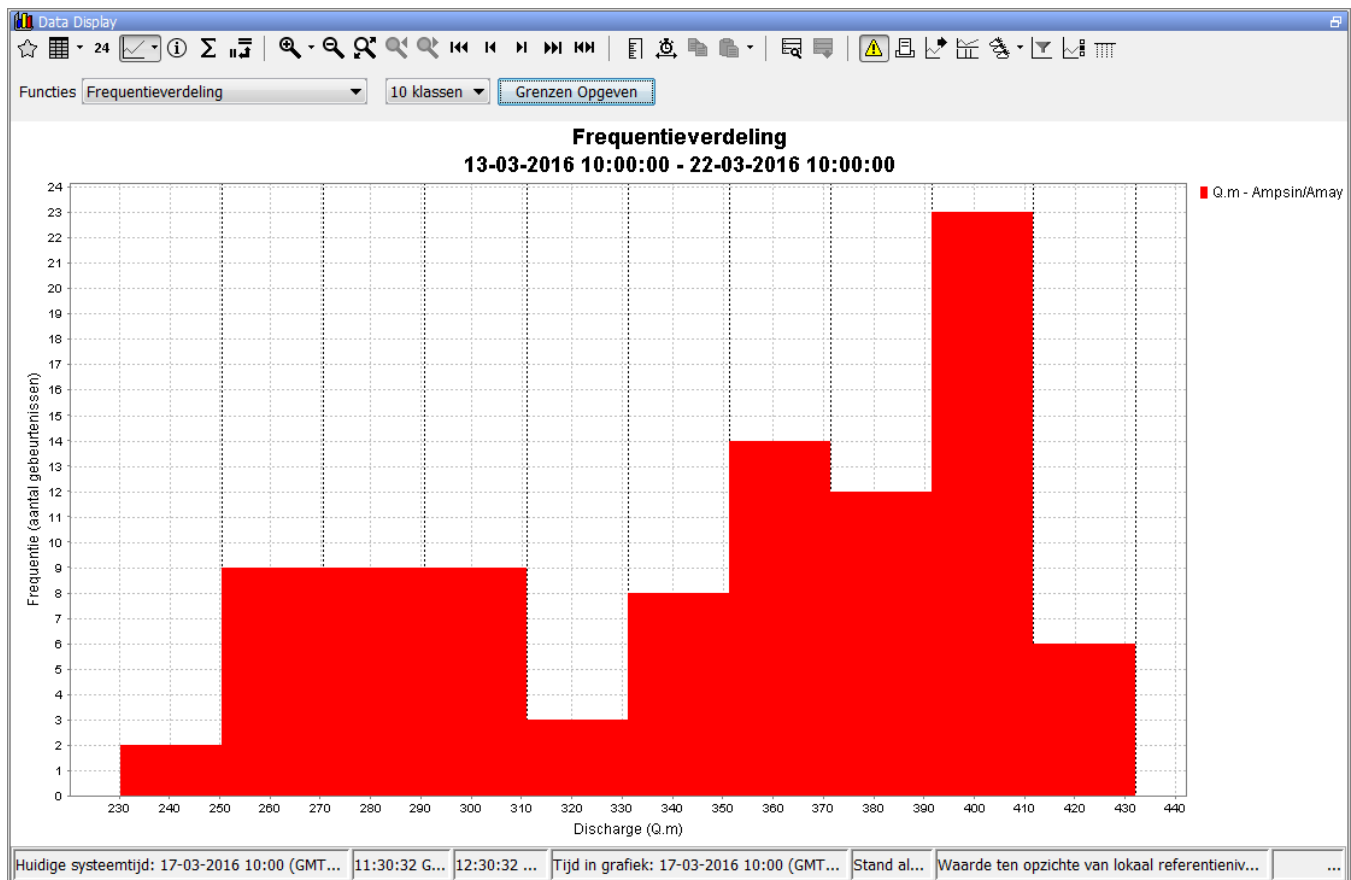


Configuration example:

```
<statisticalFunction function="ensemblePercentileExceedence"/>
```

# Frequency distribution

The frequency distribution function divides the distance between the min and max value of the timeseries by the number of samples to create a classification. It then evaluate each value in the timeseries and assigns it to a class. The result is a frequency distribution diagram listing the number of occurences per class. Number of classes (samples) can be selected in dropdown box. Also the class boundaries can be chosen in a popup window.

Popup window for selecting class boundaries.



```
<statisticalFunction function="frequencyDistribution">
   <samples amount="5"/>
   <samples amount="10"/>
   <samples amount="20"/>
 </statisticalFunction>
```
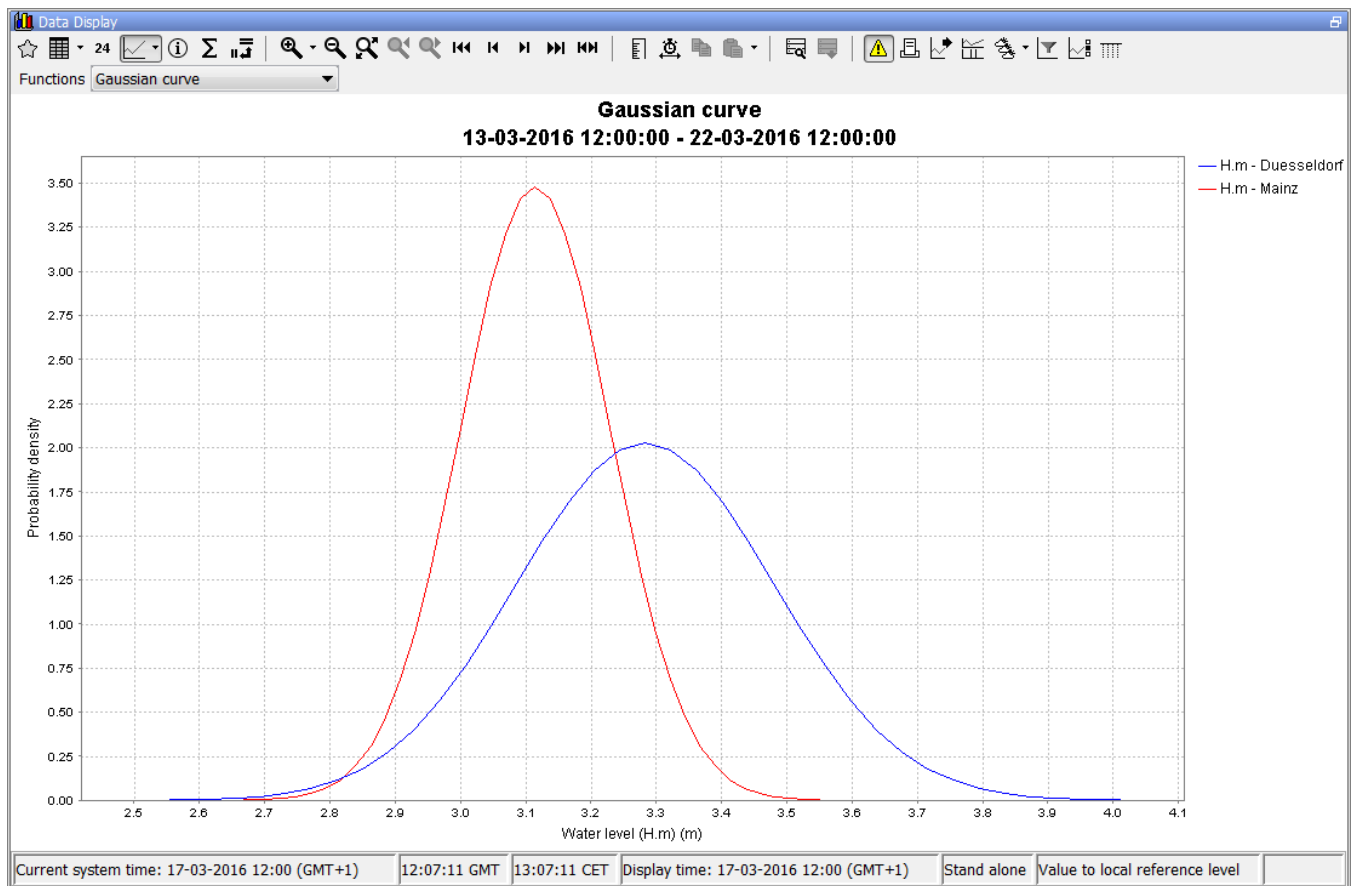
# Gaussian curve

Gaussian or bell-shaped curve showes a particular distribution of probability over the values of a random variable.

```
<statisticalFunction function="gaussianCurve" />
```

## Cumulative

Continuous accumulation over entire timeseries.



Config example:

```
<statisticalFunction function="cumulative"/>
```

## Accumulation Per Interval

Accumulation per Interval, starting at zero at the beginning of the next interval. Interval can be selected in dropdown box.



Config example:

```
<statisticalFunction function="accumulationInterval">
                    <movingAccumulationTimeSpan unit="hour" multiplier="1"/>
                    <movingAccumulationTimeSpan unit="hour" multiplier="3"/>
                    <movingAccumulationTimeSpan unit="hour" multiplier="6"/>
            </statisticalFunction>
```
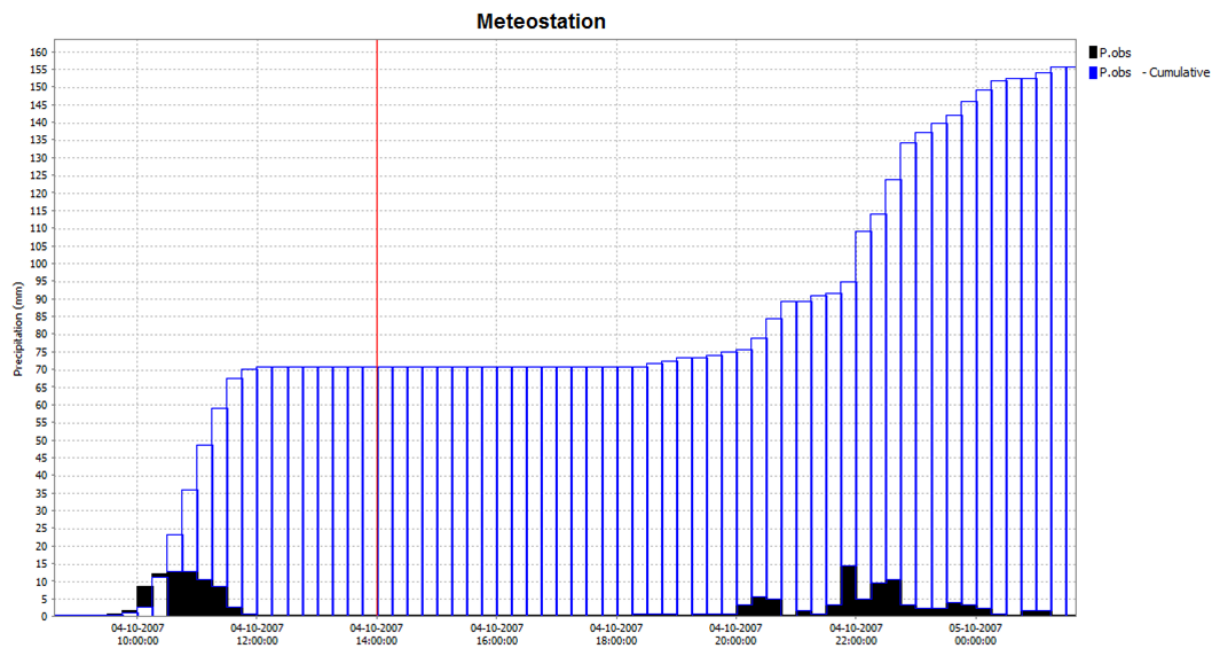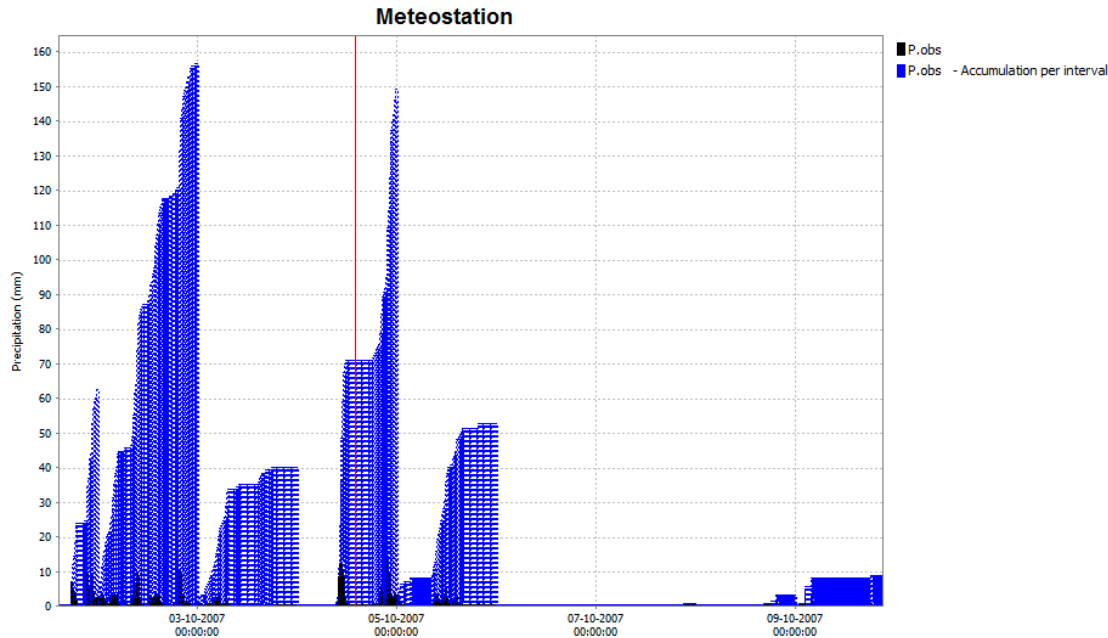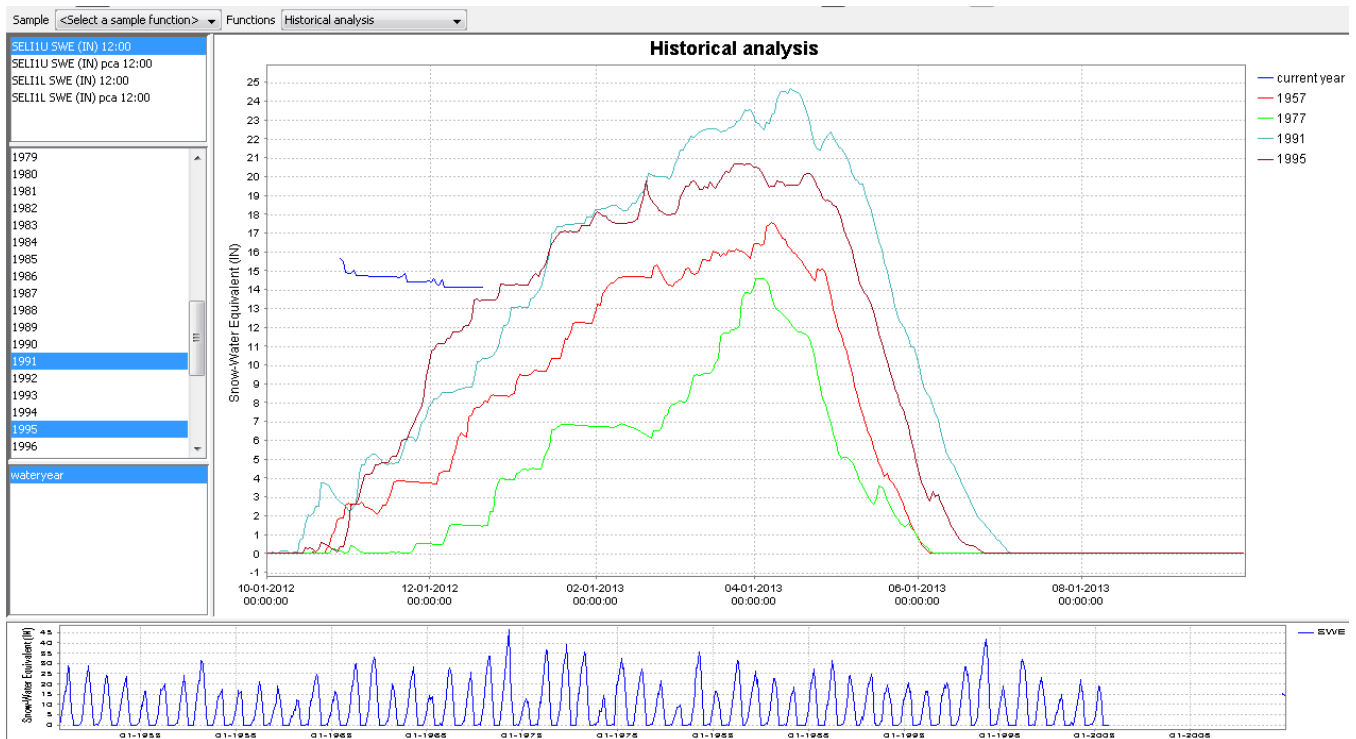
## Historical Analysis

Allows comparison of the current situation with selected previous years at the same moment within the year. Can be used to compare seasonal behaviour, e.g. deficit accumulation, snow accummulation/melt, runoff.
The bottom of the window shows the full timeseries to assist in picking relevant historical years.

The function requires a multi-year historical series, where the view period on the x-axis stretches over multiple years before the function is selected. You can use the |<>| button to stretch the x-axis from the current view period to the full available length. The display requires a fixed season definition for the x-axis, to be included in the configuration. The user needs to select the historical year of interest to plot this against the current year. Multiple years can be selected by holding the CTRL-key. Holding the SHIFT-key will select a range of years.

```
<statisticalFunction function="historicalAnalysis">
    <historicalPeriods>
        <historicalPeriod>
            <startForwardLookingPeriod>--01-01</startForwardLookingPeriod>
            <season startMonthDay="--01-01" endMonthDay="--12-01" label="wateryear"/>
        </historicalPeriod>
    </historicalPeriods>
</statisticalFunction>
```
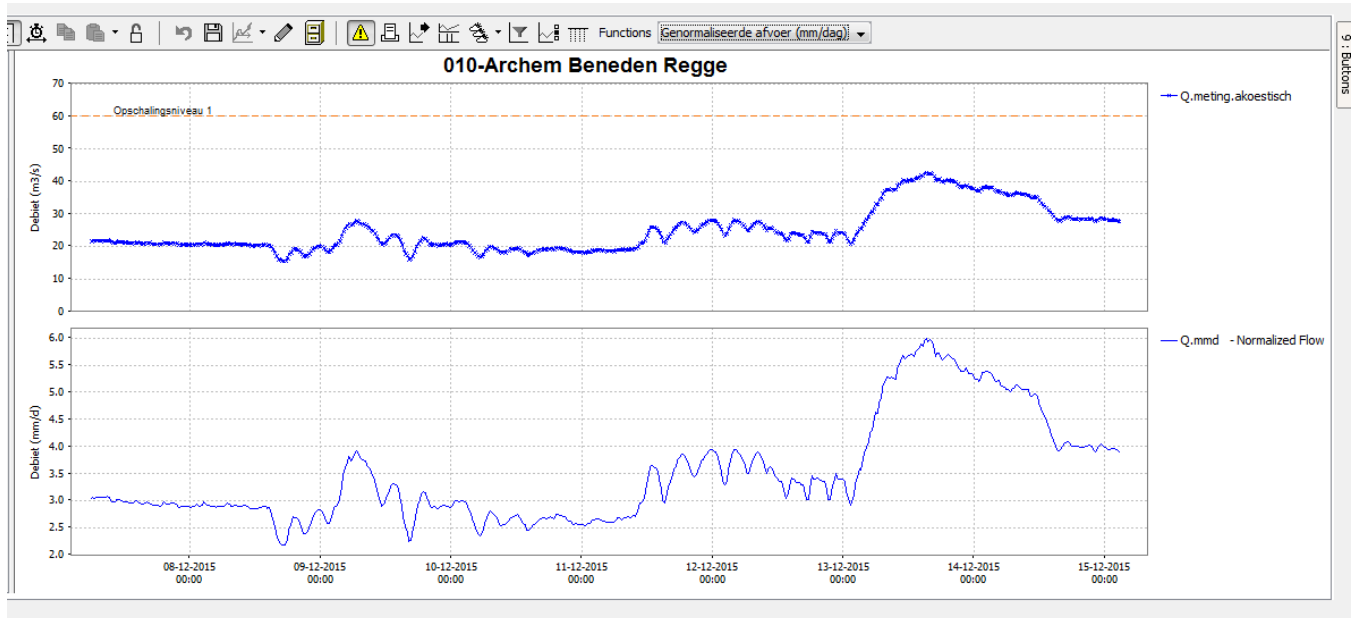
# Normalized Flow

Allows comparison of the current situation with selected previous years at the same moment within the year. CAn be used to compare seasonal behaviour, e.g. deficit accumulation, snow accummulation/melt, runoff.
The bottom of the window shows the full timeseries to assist in picking relevant historical years.

The function requires a multi-year historical series, where the view period on the x-axis streches over multiple years before the function is selected. You can use the |<>| button to stretch the x-axis from the current view period to the full available length. The display requires a fixed season definition for the x-axis, to be included in the configuration. The user needs to select the historical year of interest to plot this against the current year. Multiple years can be selected by holding the CTRL-key. Holding the SHIFT-key will select a range of years.

010-Archem Beneden Regge

```
<statisticalFunction function="normalizedFlow" label="Flow (mm/day)" ignoreMissings="true">
   <areaFunction>@AREA_HA@*10/86400</areaFunction>
   <parameterId>Q.mmd</parameterId>
   <allowedInputParameterId>Q.obs</allowedInputParameterId>
</statisticalFunction>

....and / or ....

<statisticalFunction function="normalizedFlow" label="Flow (l/s.ha)" ignoreMissings="true">
   <areaFunction>@AREA_HA@/1000</areaFunction>
   <parameterId>Q.lsha</parameterId>
   <allowedInputParameterId>Q.obs</allowedInputParameterId>
</statisticalFunction>
```

## Principal Component Analysis

The Principal Component Analysis function uses independent historical data (observations) and dependent data (e.g. a simulated basin value) to compute a number of regression equations using the Principal Component Analysis technique. The resulting equation is applied with current observations to estimate the current basin value.
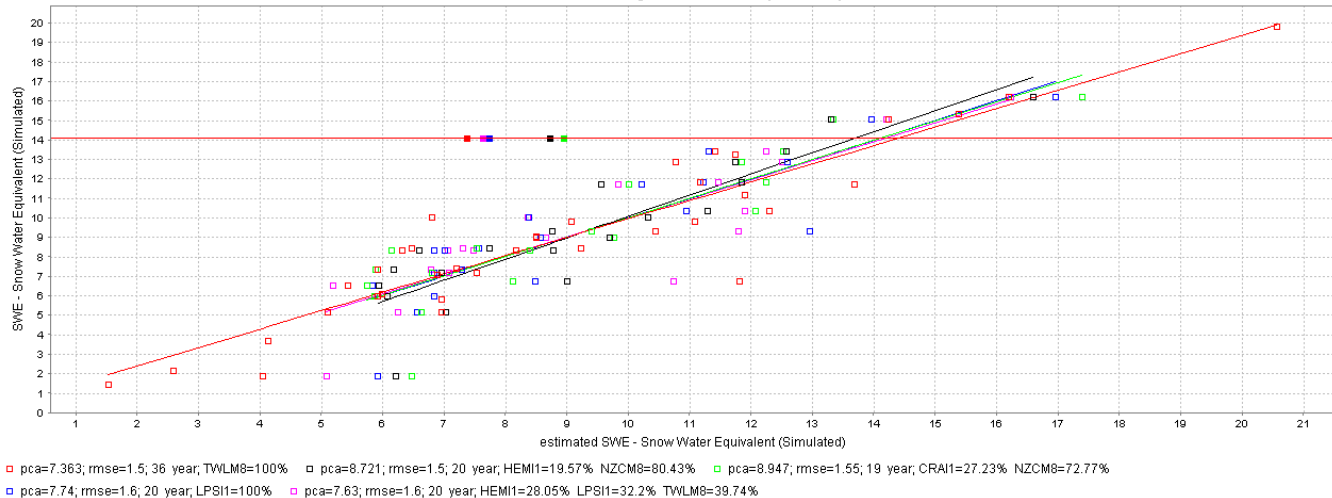
This functionality starts with the timeseries data available in the display. The independent and dependent parameters should correspond to the specification in the TimeSeriesDisplayConfig for the simulated (dependent) and observed (indenpedent) parameter.

Typically, this functionality is intended to work in combination with the Topology and the Filters, such that you can have a default set of locations which you can modify from the map or via the Filters. You can also exclude locations from the analysis by making the timeseries invisible in the graph. By default the function shows the scatterplot of the best 5 equations (lowest root mean square error). By selecting one item in the legend, this item and its confidence interval is shown.

The PCA-estimate for the basin value can be utilized in a modifier if the ModifierTypes-configuration refers to the statistical function for the default value.

| rank | simulated | pca | rmse | period of record | equation | max | min | mean | method | preprocessing | normalized |
|------|-----------|------|------|------------------|----------|------|------|------|--------|---------------|------------|
| 0 | 14.1 | 7.36 | 1.5 | 1967 - 2002 | Y = 0.809*TWLM8 -2.106 | 19.81 | 1.45 | 8.94 | PCA | NONE | no |
| 1 | 14.1 | 8.72 | 1.5 | 1983 - 2002 | Y = 0.152*HEMI1 + 1.694*NZCM8 -0.97 | 16.21 | 1.91 | 9.31 | REGRESSION | SQUAREROOT | no |
| 2 | 14.1 | 8.95 | 1.55 | 1984 - 2002 | Y = 0.207*CRAI1 + 1.707*NZCM8 -1.84 | 16.21 | 1.91 | 9.27 | REGRESSION | SQUAREROOT | no |
| 3 | 14.1 | 7.74 | 1.6 | 1983 - 2002 | Y = 0.913*LPSI1 + 0.528 | 16.21 | 1.91 | 9.31 | REGRESSION | NONE | no |
| 4 | 14.1 | 7.63 | 1.6 | 1983 - 2002 | Y = 0.18*HEMI1 + 0.328*LPSI1 + 0.273*TWLM8 -0.411 | 16.21 | 1.91 | 9.31 | PCA | NONE | no |

**PCA 12-20-2012**
**SELI1U - Selway R nr Lowell (UPPER)**



□ pca=7.363; rmse=1.5; 36 year; TWLM8=100%   □ pca=8.721; rmse=1.5; 20 year; HEMI1=19.57% NZCM8=80.43%   □ pca=8.947; rmse=1.55; 19 year; CRAI1=27.23% NZCM8=72.77%
□ pca=7.74; rmse=1.6; 20 year; LPSI1=100%   □ pca=7.63; rmse=1.6; 20 year; HEMI1=28.05% LPSI1=32.2% TWLM8=39.74%

## Scatter Plot
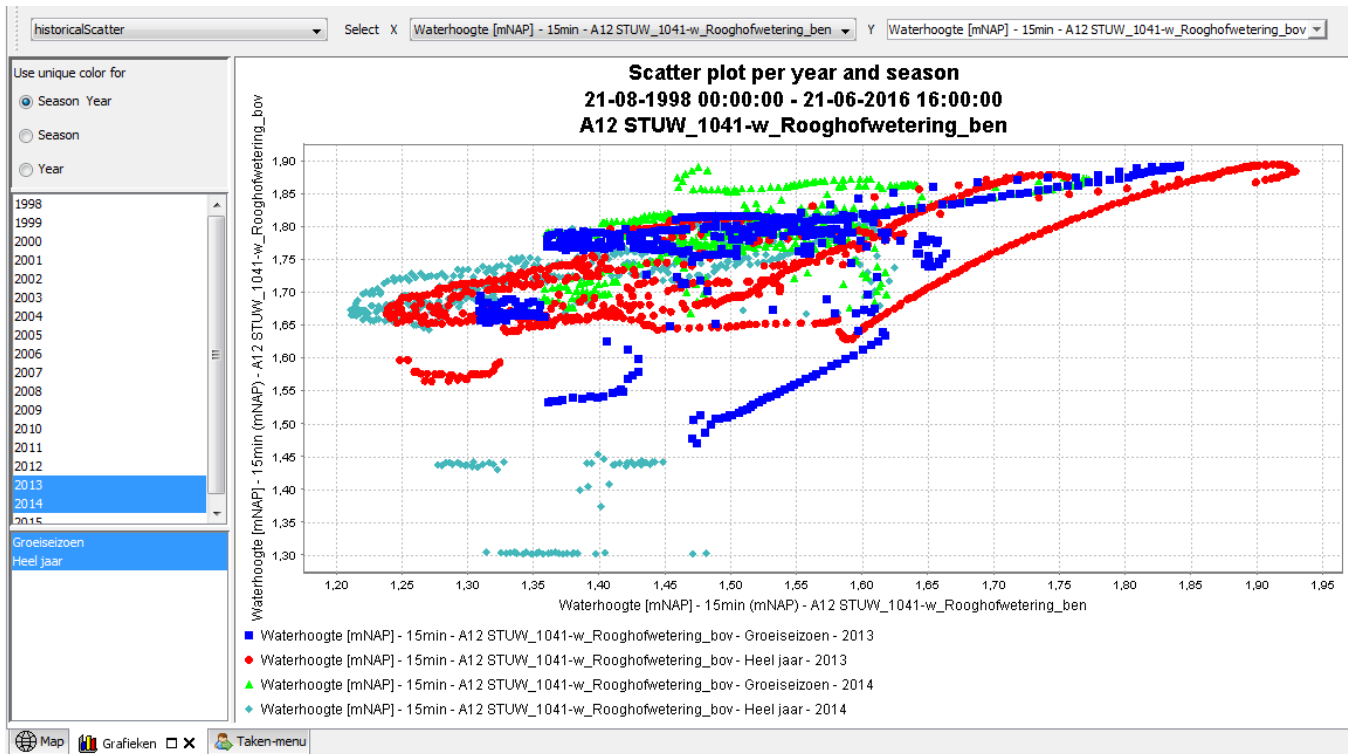
Creates a scatter plot of two selected series

Config example:

```
<statisticalFunction function="scatterPlot"/>
```

## Historical Scatter Plot (per year and season)

(Since 2016.02) There is also a scatter plot available where you can select years and seasons separately and choose unique coloring

**Example config**

```
<statisticalFunction function="historicalScatterPlot">
        <season startMonthDay="--01-01" endMonthDay="--12-31" label="Year" color="black"/>
        <season startMonthDay="--04-01" endMonthDay="--09-30" label="Grow Season" color="gray"/>
</statisticalFunction>
```

# Show statistics for specific timeseries

By default the statistical results are displayed for all timeseries. To display the statistics for a specific timeseries, the relevant timeseries can be selected by clicking the legend or by selecting the timeseries in the graph. Using the CTRL button more than one timeseries can be selected this way.

# Hide original time series

Since 2016.02 a checkbox has been added which offers the user the choice to hide the original time series. The status of this checkbox will be stored in the user settings so choice will be remembered between different statistical functions and when restarting FEWS:





The images show the checkbox for 2 functions but this checkbox is available for each statistical function which plotted the original time series together with the statistical time series.