

# **ACMANT<sub>v4</sub>: SCIENTIFIC CONTENT AND OPERATION OF THE SOFTWARE**

**Peter Domonkos**

**2<sup>nd</sup> (corrected) release**

**Tortosa, 2021**

## PREFACE

The first version of the homogenization method ACMANT was published in 2011, and that was only for the homogenization of monthly temperature series of extratropical regions. However, the method has been characterized by the accurate reconstruction of climatic trends and variability from its first version, at least when the spatial density and spatial correlations of the observed data favour statistical homogenization. The success of the first version gave the inspiration for the further developments. ACMANTv2 and ACMANTv3 were used (2014, 2015) in international trainings supported by the World Meteorological Organization. The software development was supported by the University Rovira i Virgili (Tarragona, Spain) until 2015, thereafter it has been continued from private sources.

With the development and release of newer versions, the functionality of ACMANT has been widened, and its accuracy improved further. To find the optimal homogenization for a given dataset, both the properties of the potential homogenization methods and the specifics of the homogenization task must be considered. According to present knowledge, ACMANT would likely be the most appropriate method for the homogenization of large datasets of national meteorological services, as the homogenization results of ACMANTv4 are accurate on all time scales and all spatial scales, and the use of the method is easy even for very large datasets.

The most recent version is applicable for the homogenization of several climatic variables on either monthly or daily scales. This document presents the scientific content and operation of ACMANTv4.3, while most of the technical details are presented in the Manual of the software.

I thank Victor Venema for his advice on the edition of this document.

Correction for the second release: In the first edition of this document values of parameter  $p_2$  (particular coefficient in the modified Caussinus - Lyazrhi criterion, formula 22) were 50% lowered in comparison with their true values for forgetting that in my experiments the particular coefficients were applied as additional coefficients to the standard coefficient 2.0.

Tortosa, 26-June-2021

Peter Domonkos

## Table of content

I INTRODUCTION.....	5
II BASIC CONCEPTS AND DEFINITIONS.....	7
A1 Concepts and definitions.....	7
A2 Mathematical symbols.....	10
A3 Quality indicators.....	12
III BASIC OPERATIONS.....	13
B1 Transformation of precipitation data.....	13
B2 Removal of climatic seasonality.....	13
B3 Spatial correlation.....	14
B4 Candidate series, reference series and relative time series.....	14
B5 Creation of multiple sets of reference composites.....	15
B6 Selection of relative time series.....	17
B7 Step function fitting.....	17
B8 Bivariate detection.....	18
B9 Fitting modified step function with sections including sinusoid changes.....	18
B10 ANOVA correction model.....	19
B11 Weighted ANOVA model.....	19
B12 Gap filling for additive variables.....	20
B13 Gap filling for precipitation data.....	20
IV PREPARATORY STEPS.....	22
1 Setting parameters.....	22
2 Basic control of input data.....	22
3 Construction of networks.....	23
4 Preliminary operations.....	25
5 Infilling data gaps.....	26
V FIRST ITERATION.....	29
6 Creation of relative time series for outlier filtering.....	29
7 Filtering of spatial outliers.....	30
8 Infilling data gaps.....	36
9 Creation of relative time series for break detection.....	37
10 Break detection.....	38
11 Adjustments for inhomogeneities.....	39
VI SECOND ITERATION.....	43
12 Creation of relative time series for outlier filtering.....	43
13 Outlier filtering.....	43
14 Infilling data gaps.....	43
15 Creation of relative time series for break detection.....	44
16 Break detection.....	45
17 Adjustments for inhomogeneities.....	46
VII THIRD ITERATION.....	48
18 Creation of relative time series for outlier filtering.....	48
19 Outlier filtering.....	48
20 Infilling data gaps and preliminary calculations for ensemble homogenization..	48
21 Creation of relative time series.....	50

22 Break detection.....	50
23 Adjustments for inhomogeneities.....	54
VIII FINAL OPERATIONS.....	59
24 Irregular seasonal cycle of inhomogeneities.....	59
25 Refinement of outlier periods.....	62
26 Infilling data gaps.....	64
27 Elimination of physical outliers.....	66
IX LITERATURE.....	69
L1 Description of earlier ACMANT versions and sources of ACMANT.....	69
L2 Properties of ACMANT or its routines.....	70
L3 ACMANT or its routines in method comparison studies.....	70

# I INTRODUCTION

ACMANTv4 is a software for removing non-climatic biases from climatic time series, or with the more usual term: homogenizing climatic time series. The homogenization with ACMANTv4 is based on the exploitation of the spatial redundancy existing in the observed climatic data, therefore the software needs the use of spatially sufficiently correlating networks of time series. After the preparation of input data and introduction of some parameters (e.g. number of time series, study period of time series, etc.), the operation of ACMANTv4 is fully automatic. Automatic homogenization methods have three advantages in comparison with manual or partly manual methods: i) Automatic methods can be tested in large benchmark datasets, thus their performances are the best controlled; ii) Their use is feasible even for very large data bases; iii) Their use is relatively easy. Note that the computational time demand must be kept below reasonable limits, and it is solved in ACMANTv4.

The purpose of this document is to provide the full description of the scientific content of ACMANTv4 for interested readers. No reference to other studies is given within the description, as the aim is to provide a document which is understandable without opening other documents. However, following the method description, a selection of the most relevant literature is presented.

ACMANTv4 can be applied to the homogenization of various climatic variables either on daily or monthly time scales. The execution of homogenization may differ according to climatic variables, temporal resolution of the data, presumed annual cycle of inhomogeneities and the size of the dataset. The software package includes 24 programs specific for homogenization tasks, and a high user-friendliness is provided by the inclusion of one more program which manages the homogenization procedure. The managing program performs some preparatory steps, then selects the most appropriate homogenization program of the 24 and transmits the task to the selected program.

ACMANT was constructed on the base of some earlier developed modern homogenization tools (ACMANT = Adapted Caussinus-Mestre Algorithm for the homogenization of Networks of climatic Time series), and during its development by the inclusion of new statistical tools or by other algorithm modifications its performance has been continuously tested. This empirical control is necessary, as theoretically excellent statistical tools do not always provide the optimal solution for finite and noisy samples, like observed climatic time series. According to known method comparison tests, ACMANT most often provides the smallest residual root mean square error and trend bias among the tested methods.

Documental information of changes in the performance of observations (so-called metadata) is often used together with statistical homogenization. However, ACMANTv4 cannot be used together with metadata. The utilization of metadata within automatic homogenization procedures is not a simple problem, as the reliability of the quantification of the pieces of metadata has not been profoundly studied yet.

The description of the software is organised to 7 main sections beyond to this Introduction. Section II includes a list of concepts and definitions, as well as the explication of all symbols used in the document. The concepts and definitions presented in Sec. II (24 entries) are expected to be known before the reading of the further sections. Section III presents the most important operations of ACMANTv4, many of these operations are performed repeatedly during the homogenization procedure. The remaining sections (Sec. IV – VIII) describe all the steps of the homogenization procedure in the temporal order of their execution. In Sec. IV the preparatory operations

are presented, such as the construction of climatically semi-homogeneous station networks of large datasets, calculation of derived variables from the input data, etc. Sections V – VII describe the 3 cycles of homogenization. The homogenization is always performed on the difference series of a candidate series and a composite reference series, and such difference series are named relative time series. Each homogenization cycle includes the construction of relative time series, filtering of spatial outlier values, gap filling, break detection and adjustments for the detected inhomogeneities, but the details of the execution often differ according to the specifics of the homogenization task (e.g. climatic variable, seasonality of inhomogeneities, etc.) and the phase of the homogenization procedure. With the repetition of the same kind routines, most steps are performed with gradually increasing accuracy during the homogenization procedure. The last section (Sec VIII) presents the steps which are applied after the three homogenization cycles have been finished. Some steps of this section improve more the accuracy of homogenized data (e.g. the calculation of the monthly adjustment terms for irregular shaped seasonal cycle, final gap filling, etc.), while some other steps serve the preparation of the final output results.

In spite of the easy-to-use construction of ACMANTv4, its use is recommended for skilled persons, as any error in the input data preparation or introduction of parameters may lead to serious errors in the results. The input data preparation, the properties of the output items, and some other technical details are described in the Manual of ACMANTv4.

## II BASIC CONCEPTS AND DEFINITIONS

### A1 Concepts and definitions

**Input dataset:** Collection of time series of observed climatic values. Each time series must contain the values of the same climatic element and with the same temporal resolution, which can be monthly or daily. An input dataset may have 4 to 5000 time series, which may cover varied periods.

**Network:** If an input dataset has no more than 40 time series and all of the spatial correlations are higher or equal with 0.4, then the dataset forms one only network, referred as 1-network homogenization, otherwise it will be divided to smaller networks. In ACMANTv4 the network formation is automatic (see Sect. IV/3), and when it is included, the procedure is referred to as multi-network homogenization. The number of time series within network is maximised by 99 and rarely higher than 50.

**Central series:** In multi-network homogenization each network has a central series. The other time series of the network are gathered to the network according to their spatial correlation with the central series.

**Additive variables:** In the modelling with ACMANT, the magnitude of inhomogeneities (i.e. the deviation of the observed values from the true climate) can be independent from the climatic value (additive) or proportional to the climatic value (multiplicative). The additive model is applied for temperature (**TT**), relative humidity (**HH**), sunshine duration (**SS**), radiation (**RS**), wind speed (**FF**), atmospheric pressure at station level (**PP**) and sea level pressure (**SP**). Hereafter the latter group of climatic elements is referred as additive variables.

**Multiplicative variable:** Precipitation amount (**RR**)

**Daily (monthly) homogenization:** Homogenization based on an input dataset of daily (monthly) resolution.

**Monthly value in daily homogenization:** It is monthly mean of daily values for TT, HH, FF, PP and SP always. It is monthly mean of daily values also for SS and RS in inner operations of the program, and in the output anomaly series. It is monthly total of daily values for RR always and for SS and RS in all output items except in anomaly series.

**Transformed precipitation (TR):** The multiplicative variable RR is transformed to the additive variable TR (see Sect. B1) for most of the operations with precipitation data. When operations are performed with the untransformed RR, they will be indicated in the description. Note that in spite of TR behaves as an additive variable, it does not belong to the group of additive variables in this description.

**Missing value:** Missing values are generally coded with -999.9, although the Manual describes cases when they are omitted in the input dataset. (See the rules of input dataset preparation in the Manual.)

**Missing monthly value in case of daily resolution of data:** For additive variables: when no more than 7 observed daily values are missing in a month, the status of the monthly data is observed, while it is missing or interpolated in the opposite case. For RR: when any daily data is missing, the monthly data is set to missing or interpolated.

**Missing annual value:** When no more than 3 monthly values are missing in a year, the status of annual data is observed, while it is missing or interpolated in the opposite case.

**Excluded year:** When less than 3 stations have observed annual values for a year, that year is excluded from most steps of the homogenization procedure, exceptions will be indicated.

**Target period of homogenization:** While the periods of input time series may be varied, the target period is fixed for the entire dataset. A target period can be either shorter or longer than the period covered by individual time series of the input dataset. When an input time series is longer than the target period, the input data outside the target period will be dropped, and when an input time series is shorter than the target period, that series will be supplied automatically with missing data codes before the homogenization. The target period is defined by the user at the beginning of the homogenization. Its length is expected to be between 10 and 200 years, and ideally it reflects well both the purpose of homogenization and the data availability. After the definition of the target period, all the time series of the dataset have the same length. An example of  $N$  monthly series of  $n$  year period is shown by (1).

$$\mathbf{X}_s = x_{s,1}, x_{s,2}, \dots, x_{s,h} \dots x_{s,H} \quad (s = 1, 2 \dots N), \quad H = 12n \quad (1)$$

In (1),  $h$  stands for the serial number of month from the beginning of the time series.

**Treated period of time series:** ACMANT needs a certain amount and temporal compactness of the observed data for treating time series. The minimum amount of observed monthly values is 114 and the minimum compactness is 25% temporal density of observed values. Note, however, that between two adequately compact blocks with at least 60 observed monthly values in each, the extent of data gaps is unlimited.

A treated period includes entire years only, i.e. it starts with 01 January of its first year and ends with 31 December of its last year. The minimum length of treated period is 10 years without excluded years. Time series without acceptable treated period are left out of consideration during the homogenization procedure. See further conditions for treated periods in the Manual. Usually no or few observed data occur outside the treated period. Observed data outside the treated period are left out of consideration in most steps of the homogenization procedure, exceptions will be indicated.

**Homogenized period of time series:** ACMANT needs at least 4 time series of sufficiently high spatial correlations for performing homogenization, and this condition is also a requirement for sections of time series. The homogenized period includes entire years only, and it is shorter than or identical to the treated period. Time series without homogenized period are left out of consideration in most steps of the homogenization procedure.

Note that this concept reflects the spatial-temporal coherence of the data, and does not the status of the data, so that it can be applied either to series already



homogenized or to series those will be homogenized in a later phase of the homogenization procedure.

If a given year would be part of the homogenized periods of less than 4 station series of a network, this year will not be part of a homogenized period in any series. If more than one homogenized periods appear for a station series, only the last of them will be used as homogenized period.

**Station effect:** Non-climatic component ( $v$ ) of observed data, influenced by the characteristics of station location, as well as the technical and personal conditions of climate observation.

$$\mathbf{X}_s = \mathbf{U}_s + \mathbf{V}_s + \boldsymbol{\varepsilon}_s \quad (2)$$

In (2)  $\mathbf{U}$  stands for the regional climate signal,  $\mathbf{V}$  for the station effect, while in  $\boldsymbol{\varepsilon}$  the effects of weather and occasional observation errors are summed up and the subscript  $s$  is the station index.  $\boldsymbol{\varepsilon}$  usually can be modelled well with white noise or red noise. In case of a homogeneous series, the station effect is constant, while temporal changes of  $v$  are inhomogeneities.

**Break:** Sudden shift in the section means of station effect. A break can be characterised with its date and size. The model assumes a break is instantaneous; it takes places between two consecutive values of time series. The date of the break is the last date before the event.

**Outlier value:** Two kinds of outlier values are considered: physical outliers and spatial outliers. The thresholds for physical outliers can be user-defined values or default values (see Manual). In monthly homogenization, spatial monthly outliers of the homogenized period are filtered (optional). Possible spatial daily outliers are not examined.

**Outlier period:** Short-term biases, examined in the homogenization of additive variables. These inhomogeneities can be detected only for relatively large size biases, and in this respect they are similar to spatial outliers. In ACMANTv4 the length of outlier periods varies between 10 days and 28 months (2 months and 28 months) in daily (monthly) homogenization.

**Seasonal cycle:** The modeled seasonal cycle of station effects (user defined, see Manual) is referred with this term, it can be sinusoid, irregular or constant.

**Summer – winter difference:** When the seasonal cycle is sinusoid, the summer – winter difference (denoted with upper wave) is used in the homogenization procedure. It is defined for any month  $h_0$  of  $\mathbf{X}$  by the values within a time window around  $h_0$  (3).

$$\widetilde{x}_{h_0} = \frac{1}{3.5} (\sum_{h=h_0-5}^{h_0+5} \mu_m x_h + 0.5 \mu_m (x_{h_0-6} + x_{h_0+6})) \quad (3)$$

For summer months:  $\mu_5 = \mu_6 = \mu_7 = 1$        $\mu_8 = 0.5$   
 For winter months:  $\mu_1 = \mu_{11} = \mu_{12} = -1$        $\mu_2 = -0.5$   
 For the other months:  $\mu_3 = \mu_4 = \mu_9 = \mu_{10} = 0$

When  $h_0$  is closer than 6 months to the starting or ending month of  $\mathbf{X}$ , the time window around  $h_0$  will be truncated and a logically fitting definition is provided (not shown). The relation between the serial number of month from the beginning of time series  $h$ , the serial number of year ( $y$ ) from the beginning of time series and the within year serial number of calendar month  $m$  is defined by (4).

$$h = 12(y - 1) + m \quad y \in \{1, 2, \dots, n\} \quad (4)$$

The definition of summer – winter difference for year  $y$ , coherent with (3,4), is shown by (5).

$$\widetilde{x}_y = \frac{1}{3.5} (\sum_{m=1}^{12} \mu_m x_{y,m}) \quad (5)$$

**3-month overlapping seasons:** 12 seasons constructed by merging 3 adjacent calendar month, i.e. JFM, FMA, MAM, etc. They will be referred to as 3-month seasons.

**Bi-seasons for RR:** Rainy season and snowy season. It must be uniform for a dataset. A month belongs to the snowy season if more than 50% of the precipitation falls in form of solid precipitation, while it belongs to the rainy season in the reverse case. A seasonal value is the sum of the monthly values. Under a warm climate the full year is considered the rainy season (1-season model). Note that “dry season” does not exist in ACMANT homogenization. The length of the snowy season can be 0 or between 3-9 months. If it was 1 month or 2 months by user definition, the adjacent months are put together with the snowy months to lengthen the snowy season. If it was higher than 9 months by user definition, the 1-season model is applied.

**Ensemble homogenization:** Some steps of homogenization are repeated with slightly differing setups for the ensemble members and statistics of the ensemble results are used for adjusting inhomogeneities. Steps of the homogenization procedure repeated within an ensemble cycle will be marked with “E” in the headline of the step.

## A2 Mathematical symbols

Most of the symbols, although not all of them, are used with the same meaning throughout this document. The meanings of letters  $a$ ,  $A$ ,  $b$ ,  $B$ ,  $i$  and  $j$  are varied and they are defined at the relevant section of the document. Versions of the same kind variable are often distinguished with apostrophe, asterisk or other supplements. The list of symbols shown here always includes the basic version of the symbols, and only for a few cases includes other versions. Most of the mathematical symbols are printed with italics, except for  $Y$  (which represents cluster) and vectors of time series, the latter ones are printed bold.

$d$  – calendar day

$d'$  – serial number of the day from the beginning of the examined section of time series

$D$  – number of days within a month

$E$  – external variance

$f, \mathbf{F}$  – reference series  
 $g, \mathbf{G}$  – deseasonalised series  
 $gc, \mathbf{Gc}$  – candidate series (deseasonalised)  
 $h$  – serial number of the month from the beginning of the examined section of time series  
 $H$  – number of months in time series  
 $I$  – internal variance  
 $J$  – number of effective partners  
 $k$  – serial number of break  
 $K$  – number of breaks in a time series  
 $l$  – length of period in days  
 $L$  – length of period in months  
 $m$  – calendar month  
 $m^*$  – 3-month season  
 $M$  – number of relative time series for a given candidate series  
 $n$  – number of years in the treated period  
 $n'$  – number of years in the homogenized period  
 $N$  – number of stations in the dataset  
 $N'$  – number of stations in a network  
 $N^*$  – number of reference composite series  
 $p$  – parameter  
 $P$  – penalty term  
 $q$  – bias size for a section of time series  
 $Q$  – score for break magnitude  
 $r$  – spatial correlation  
 $r^*$  – spatial correlation of increment series  
 $s$  – station index of time series  
 $S$  – score  
 $t, \mathbf{T}$  – relative time series  
 $u, \mathbf{U}$  – climate signal  
 $v, \mathbf{V}$  – station effect  
 $w$  – weight  
 $W$  – sum of the weights of the partner series  
 $x, \mathbf{X}$  – climate data series  
 $xc, \mathbf{Xc}$  – candidate series  
 $\tilde{x}, \tilde{\mathbf{X}}$  – summer – winter difference  
 $y$  – serial number of year from the beginning of the examined section of time series  
 $Y$  – cluster  
 $z, \mathbf{Z}$  – adjustment term  
 $\alpha$  – statistical significance  
 $\beta$  – usefulness of relative time series  
 $\delta$  – break size  
 $\Delta$  – deviation  
 $\varepsilon$  – noise  
 $\lambda$  – distance in years  
 $\mu$  – season-coefficient  
 $\sigma$  – standard deviation  
 $\omega, \mathbf{\Omega}$  – adjustment term for the seasonal variation of biases

### A3 Quality indicators

Series **G**, **Gc**, **X**, **Xc** often hold quality indicators, which show that the time series have passed the outlier filtering and/or inhomogeneity adjustments, or not.

**G** – neither outlier filtered, nor homogenized

**G<sup>+</sup>** – outlier filtered, but not homogenized (note that outlier filtering includes the removals both of single monthly outlier values and the values of outlier periods when the procedure includes these steps)

**G<sup>#</sup>** – outlier filtered for outlier periods shorter than 5 months. Longer outlier periods and inhomogeneities have not been treated.

**G<sup>\*</sup>** – pre-homogenized, but not outlier filtered

**G<sup>+\*</sup>** – outlier filtered and pre-homogenized

**G<sup>\*-</sup>** – not outlier filtered, but completely homogenized except for the seasonal variation of biases.

**G<sup>\*\*</sup>** – outlier filtered and completely homogenized, except for the seasonal variation of biases.

**G<sup>\*\*</sup>** – completely homogenized and outlier filtered

### III BASIC OPERATIONS

Operations which are both important and executed in the same way in different stages of the homogenization procedure are presented here. Details dependent on the stage of the homogenization procedure are not shown in this section.

#### B1 Transformation of precipitation data

A quasi-logarithmic transformation is applied to RR data to create the additive version (TR) of this variable (6).

$$\begin{aligned} TR &= \ln(RR) && \text{when } RR \geq 30 \text{ mm} \\ TR &= \ln(0.4RR + 0.01RR^2 + 9) && \text{when } RR < 30 \text{ mm} \end{aligned} \quad (6)$$

Monthly, annual and bi-seasonal RR are subjected to this transformation. Annual and bi-seasonal values of TR are converted from annual and bi-seasonal RR, respectively, and never from summing up monthly TR.

#### B2 Removal of climatic seasonality

For additive variables, as well as for monthly TR: As a preparatory step, the seasonal cycle is removed and it is added back only at the end of the homogenization procedure. The treatment for TR is partly different (see step 4.5).

Observed values of the treated period of series  $s$  are separated, and their cluster is denoted with  $Y_s$ , its sub-cluster for calendar month  $m$  with  $Y_{s,m}$ . The number of elements in  $Y_s$  and  $Y_{s,m}$  are  $H_s^*$  and  $H_{s,m}^*$ , respectively. The monthly climatic normal ( $\overline{X_{s,m}}$ ) is subtracted from the observed values (7-9).

$$\overline{X_{s,m}} = \frac{1}{H_{s,m}^*} \sum_{Y_{s,m}} x_{s,h} \quad (7)$$

$$g_{s,h} = x_{s,h} - \overline{X_{s,m}} \quad (8)$$

$$\mathbf{G}_s = g_{s,1}, g_{s,2}, \dots, g_{s,h} \dots g_{s,H} \quad (9)$$

Upper stroke denotes temporal average (arithmetical average). Note that (i) in daily homogenization (7) is unchanged, and in (8)  $\overline{X_{s,m}}$  is subtracted from the daily values of  $\mathbf{X}$ ; (ii) in the daily homogenization of SS (RS),  $\overline{X_{s,m}}$  represents the monthly mean climatic value of daily SS (RS).

### B3 Spatial correlation

Two kinds of spatial correlations are used, both of them are based on deseasonalised monthly values. In precipitation homogenization, the transformed variable (TR) is used. In the first version ( $r$ ) the Spearman correlation between simultaneous monthly values is calculated. In the other version ( $r^*$ ) the increment series of monthly values are calculated, and then the Spearman correlation is computed from them.  $r^*$  is widely used in time series homogenization, since it is less affected by inhomogeneities than  $r$ . In ACMANTv4 mostly  $r^*$  is used, except for gap filling, because there data of limited time windows are used, hence the impact of possible inhomogeneities is different in gap filling than in the other steps of the homogenization procedure.

Both for  $r$  and  $r^*$ : Both values of a source data pair must be observed values (i.e. interpolated values are excluded), and they must be within the treated period. The minimum number of source data pairs is 50, otherwise the correlation is considered 0 (except at step 3.5 correlations for such cases remain undetermined). Correlations below 0.4 are treated as 0. In calculating  $r^*$ , the use of the source data of the homogenized period is preferred (except before the first homogenization cycle), and only when the number of source data pairs within homogenized period would be less than 100, the period of source data pairs is extended to the treated period.

Even when sections of time series are examined, the spatial correlations are always calculated from the data of the entire treated period, in order to reduce estimation errors which would be larger for short sections.

### B4 Candidate series, reference series and relative time series

In break detection and outlier filtering steps the difference of a candidate series ( $\mathbf{G_c}$ ) and its reference series ( $\mathbf{F}(g_c)$ ) composed from other series of the same network examined and it is named relative time series ( $\mathbf{T}$ ) in this study (10-11).

$$\mathbf{T} = \mathbf{G_c} - \mathbf{F} \quad (10)$$

$$\mathbf{F} = \frac{\sum_{s=1}^{N^*} w_s \mathbf{G_s}}{\sum_{s=1}^{N^*} w_s} \quad (11)$$

$N^*$  is the number of partner series used to compute the composite reference series ( $3 \leq N^* \leq N' - 1$ ).  $\mathbf{T}$  does not contain regional climate signal, but only the local deviations from that ( $\Delta\mathbf{U}$ ), (12):

$$\mathbf{T} = \Delta\mathbf{U}_{g_c} - \Delta\mathbf{U}_F + \mathbf{V}_{g_c} - \mathbf{V}_F + \boldsymbol{\varepsilon}_T \quad (12)$$

As  $\Delta u$  are usually small ( $\Delta u \ll \varepsilon$ ),  $\mathbf{T}$  often can be modelled as a composition of the station effect of the candidate series plus a white noise or red noise. However, inhomogeneity of the reference series ( $\mathbf{V}_F$ ) might be of significant magnitude, affecting the detection of  $\mathbf{V}_{g_c}$ . Eqs. (10-12) describe the general problem of the relative homogenization.

In ACMANTv4 the possible effect of  $\mathbf{V_F}$  is attenuated with the following tools:  
a) Three iteration cycles; b) Ensemble homogenization; c) Strict rules for the composition of reference series, which are as follows.

- i) Each reference composite must have sufficient spatial correlation with the candidate series. The minimum threshold  $r_{min}^* = 0.4$ , except when  $N^* = 3$ ,  $r_{min}^* = 0.5$ .
- ii) The homogenized period of a reference composite cannot be shorter than the examined period of the candidate series.
- iii) Possible data gaps in the homogenized period of a reference composite are infilled with spatial interpolation before its aggregation to the reference series,
- iv) The common period of the candidate series and the homogenized period of the reference composites cannot be shorter than 10 years, and the minimum number of observed monthly values within this common period is 12 for each reference composite (while the other values may be either observed or interpolated).
- v) Generally, the higher number of reference composites meeting with conditions (i-iv) are used, the better results are expected, therefore the number of reference composites is unlimited except when the data are examined in higher than annual resolution.

Often, higher number of reference composites can be found for a limited section of the candidate series than to the whole homogenized period of that. Note that (10) can be applied also for selected periods of the candidate series, hence different periods may be homogenized with different reference series. In this way, relatively short reference composites can also be exploited in the homogenization procedure (see B5).

## B5 Creation of multiple sets of reference composites

For a given candidate series  $\mathbf{G_c}$ , the number of possible reference composites may vary from year to year, but within a given year it must be constant, as homogenized periods include entire years. The concept of “best fitting reference series” to year  $y$  of the candidate series will mean the composition of the maximum possible number of reference composites including year  $y$ , and its length will be the maximum period allowed by the composites. A simplification is that the term “reference series” will be used here in Section B5 for the composition of reference composites (with symbol  $\mathbf{F}$ ), leaving out of consideration that the final reference series depend on the weighting of reference composites (11). The rules of weighting vary according to the steps of the homogenization procedure and will be detailed later.

Often, the same set of reference composites is applicable for many years, hence the number of reference series is often lower than 10, in spite of their theoretical maximum number is as large as  $\frac{(n'-10)^2}{2} \approx 18,000$  when  $n' = 200$ . A procedure is established (B5.1 and B5.2) to exclude less effective versions, and to limit the maximum number of reference series by 80.

### B5.1 Covering all years of the homogenized period with at least 1 reference series

The purpose of this step is to cover all years of the candidate series with at least 1 reference series whenever it is possible. The best fitting reference series to the first year is selected, its period is  $[1, y_1]$ , and the sum of squared spatial correlations with the candidate series (13) is retained.

$$W_i = \sum_{Y_i} r^{*2} \quad (13)$$

In (13),  $Y_i$  denotes the cluster of the time series composing the reference series  $\mathbf{F}_i$ , while  $W_i$  is the sum of the weights. As a continuation, the best fitting reference series to year  $y_{i+1}$  is selected, etc. until year  $n'$  will be covered by a reference series. As the minimum length of a reference series is 10 years,  $M' \leq n'/10$  reference series are constructed and retained at this phase.

#### B5.2 Selecting the best fitting reference series for each year

Step-by-step, the best fitting reference series are selected for years  $2 \dots n'$ . In a given step, reference series  $\mathbf{F}_i'$  is compared with each of the  $M''$  reference series have already been retained ( $M'' \geq M'$ ). Let  $y_{i,1}$  ( $y_{i,n_i}$ ) stand for the first (last) year of  $\mathbf{F}_i'$ .  $\mathbf{F}_i'$  will be retained if at least one of relations (14-16) (but not necessarily the same relation) is true for all of the earlier selected reference series ( $j = 1, 2, \dots, M''$ ), while it is dropped in the reverse case.

$$0.95W_i > W_j \quad (14)$$

$$y_{i,1} \leq y_{j,1} - p_1 \quad (15)$$

$$y_{i,n_i} \geq y_{j,n_j} + p_1 \quad (16)$$

Usually  $p_1 = 0$ , but in the very unlikely case of  $M'' > 80$ , the procedure of B5.2 is re-started with a new  $p_1$  elevated with 1 in comparison with its previous value.

Finally, with fixing the weights of reference composites, the set of  $M$  reference series ( $\mathbf{F}_1, \mathbf{F}_2, \dots, \mathbf{F}_M$ ) and  $M$  relative time series ( $\mathbf{T}_1, \mathbf{T}_2, \dots, \mathbf{T}_M$ ) are constructed by Eqs. (10-11) for a given candidate series.

#### B5.3 Quality level of candidate series

ACMANTv4 includes 3 iteration cycles, and in later iterations the previous operations are usually repeated with data of increased quality. However, the quality level of candidate series in break detection and outlier filtering is exception, as the same inhomogeneities and outliers are examined in the later iterations as in the earlier ones. In break detection, the candidate series is outlier filtered, but not inhomogeneity adjusted ( $\mathbf{Gc}^+$ ). In outlier filtering, the candidate series includes the adjustments for inhomogeneities, but excludes any correction for outliers ( $\mathbf{Gc}^*$ ). By contrast, once the program has proceeded the first iteration cycle, the reference composites are both outlier filtered and homogenized ( $\mathbf{G}^{+*}$ ).



## B6 Selection of relative time series

A particular year of **Gc** may belong to several relative time series, but the detection of a break or an outlier is based on one only **T**. In the selection of **T**, higher  $W$  and longer series are preferred via indicator  $\beta$  (17).

$$\beta_i = W_i \ln(6n_i) \quad i \in (1, 2, \dots, M) \quad (17)$$

In (17), symbols are used with the same meanings as in B5. The relative time series are ordered according to their  $\beta$  scores. At any date of the candidate series, the **T** of the highest  $\beta$  will be selected from those which covering the given date. Note, however, that for reducing tail-effects in break detection close to an endpoint of **T**, sometimes not the **T** of highest  $\beta$  is used, but another relative time series in which a break or an outlier period is further from the endpoints of the series. In connection with this problem, a certain overlapping of the relative time series is applied in some specific cases. Such details will be described later.

Note that during the homogenization procedure any series of a network can be candidate series, thus indexing the relative time series both according to the candidate series and the serial number of reference series would be the most accurate. However, most operations are performed separately and in the same way for each candidate series, and with the best fitting relative time series, therefore these identifiers will usually be omitted.

## B7 Step function fitting

A step function with  $K$  steps splits the time series to  $K+1$  constant sections. This model approaches well the true properties of a relative time series including  $K$  breaks and little noise, but note that the model functions adequately even when there is a difference in the number of small breaks or when not all the inhomogeneities are breaks. We introduce the concepts of internal distance ( $I$ ) for values from their section mean and external distance ( $E$ ) for section means from the mean of the whole time series (18, 19).

$$I_y = t_y - \overline{T_k} \quad k = 0, 1, 2, \dots, K, \quad y \in k \quad (18)$$

$$E_k = \overline{T_k} - \overline{T} \quad (19)$$

For a given number of  $K$ , the step positions with minimum internal variance provides the best fitting step function (20), known also as optimal segmentation.

$$\text{Optimal solution} \equiv \min_{y_1, y_2, \dots, y_K} \sum_{k=0}^K \sum_{y=y_k+1}^{y_{k+1}} I_y^2 \quad (20)$$

In (20),  $y_1, y_2, \dots, y_K$  are the positions of breaks,  $y_0 = 0$ ,  $y_{K+1} = n$ .

In time series homogenization the optimal number of steps is unknown. Here, a modified version of the Caussinus – Lyazrhi criterion (C-L criterion) is used for setting  $K$ . In the C-L criterion the first expression monotonously decreases with increasing  $K$ ,

but this decrease may be balanced or overbalanced by the penalty term  $P$  which increases with increasing  $K$  (21,22).

$$\text{Optimal solution} \equiv \min_K \left\{ \ln \left( 1 - \frac{\sum_{k=0}^K (y_{k+1} - y_k) E_k^2}{\sum_{y=1}^{n'} (t_y - \bar{T})^2} \right) + P \right\} \quad (21)$$

$$P = \frac{p_2 K}{n' - 1} \ln(n') \quad (22)$$

Eqs. (19-22) show the example of annual resolution for the homogenized period ( $n'$  years), but note that these formulas can be applied for any time period and with any time resolution. Eqs. (21-22) differ only in parameter  $p_2$  from the original C-L criterion where  $p_2 = 2$ .

### B8 Bivariate detection

The step function fitting with minimising the sum of internal distances for two variables ( $\mathbf{T}_A$  and  $\mathbf{T}_B$ ) is shown by (23-24).

$$\min_{y_1, y_2 \dots y_K} \left\{ \sum_{k=0}^K \sum_{y=y_{k+1}}^{y_{k+1}} (I_{A,y}^2 + p_3 I_{B,y}^2) \right\} \quad (23)$$

$$\min_K \left\{ \ln \left( 1 - \frac{\sum_{k=0}^K (y_{k+1} - y_k) (E_{A,k}^2 + p_3 E_{B,k}^2)}{\sum_{y=1}^{n'} ((t_{A,y} - \bar{T}_A)^2 + p_3 (t_{B,y} - \bar{T}_B)^2)} \right) + P \right\} \quad (24)$$

### B9 Fitting modified step function with sections including sinusoid changes

When the model of inhomogeneities is sinusoid, the unbiased values are approximated by a modified step function including sections with no change in annual values, but of sinusoid seasonal changes for monthly and daily values. For climatological reasons, the modes of the oscillation are at the solstices. The transformation of (20) to (25, 26) provides the solution for modified step function fitting to data of monthly resolution.

$$\text{Optimal solution} \equiv \min_{h_1, h_2 \dots h_K, a, b} \sum_{k=0}^K \sum_{h=h_{k+1}}^{h_{k+1}} I_h'^2 \quad (25)$$

$$I_h' = t_h - a_k - b_k \sin\left(\frac{2\pi(m-3.2)}{12}\right) \quad k = 0, 1, 2, \dots K, \quad h \in k \quad (26)$$

Note that when modified step function fitting is applied in ACMANTv4,  $K \leq 2$  by definition. When a section between two steps is shorter than 10 months, the function is constant there ( $b_k = 0$ ).

## B10 ANOVA correction model

ACMANTv4 applies two versions of the ANOVA correction models, one for spatially homogeneous climate, referred to as “ANOVA correction model”, and a more complex model including the spatial variation of climate, referred to as “weighted ANOVA model”.

The ANOVA correction model is based on the model of Eq. (2), and on the minimization of the variance of homogenized time series. In the practical solution, Eq. (2) is applied on the estimates of  $\mathbf{U}$  and  $\mathbf{V}$  ( $\widehat{\mathbf{U}}$  and  $\widehat{\mathbf{V}}$ ), with known break positions and under the condition  $\hat{\varepsilon} \equiv 0$ . In this model,  $\Delta \mathbf{U} \equiv 0$  within a given network. Based on these conditions, the equation system (27, 28) is constructed, from which the optimal estimates of  $\mathbf{U}$  and  $\mathbf{V}$  can be calculated.

$$N_h \widehat{u}_h + \sum_{s=1}^{N_h} \widehat{v}_{s,k} = \sum_{s=1}^{N_h} x_{s,h} \quad \text{for every } h \in [1, H'], \quad H' = 12n', \quad k = k(s, h) \quad (27)$$

$$\sum_{h=h_k+1}^{h_{k+1}} \widehat{u}_h + (h_{k+1} - h_k) \widehat{v}_{s,k} = \sum_{h=h_k+1}^{h_{k+1}} x_{s,h} \quad \text{for every } s \text{ and } k \quad (28)$$

In (27),  $N_h$  stands for the number of stations for which  $h$  falls within the homogenized period. Eqs. (27-28) show the model for monthly resolution, but it can be applied for any time resolution.

The input data of the model is the observed time series ( $\mathbf{X}$ ) and the estimated break positions. Note, however, that the following operations are made in ACMANTv4, before inputting the data into the model: a) Climatic monthly mean vales are removed (although it impacts  $\widehat{\mathbf{U}}$ , neutral to  $\widehat{\mathbf{V}}$ ); b) Outlier filtering (it increases accuracy); c) Precipitation data are transformed to an additive variable (TR), and it is necessary to the application of the model in the present form, as inhomogeneities are expected to result in linear biases.

## B11 Weighted ANOVA model

One station series of a network is selected to be candidate series ( $\mathbf{Gc}$ ), and the other series are weighted with the squared spatial correlations.

$$\sum_{s=1}^{N_h} r_{gc,s}^{*2} \widehat{u}_{gc,h} + \sum_{s=1}^{N_h} r_{gc,s}^{*2} \widehat{v}_{s,k} = \sum_{s=1}^{N_h} r_{gc,s}^{*2} x_{s,h} \quad (29)$$

$$\sum_{h=h_k+1}^{h_{k+1}} \widehat{u}_{gc,h} + (h_{k+1} - h_k) \widehat{v}_{s,k} = \sum_{h=h_k+1}^{h_{k+1}} x_{s,h} \quad (30)$$

Eqs. (29-30) give the best solution only for  $\mathbf{Gc}$ , therefore the calculations must be repeated for each series of the network. As a consequence, the computational time demand of weighted ANOVA is  $N$  times higher than that of the “simple” ANOVA. Note that ordinary kriging would serve the optimal weights, but the inclusion of kriging would even more elevate the computational time demand, while the ANOVA weights have generally little impact on the accuracy of the results.

Regarding the specifics of the model input data preparation in ACMANTv4, they are the same as those for B10.

### B12 Gap filling for additive variables

A spatial interpolation for the monthly missing value  $h0$  of candidate series is performed using the weighted average of neighbouring observed values tuned to the long-term mean of  $\mathbf{Gc}$  ( $g'_{s,h0}$ ), (31-34).

$$gc_{h0} = \frac{1}{W'} \sum_{s=1}^{N''} w_s g'_{s,h0} \quad (31)$$

$$g'_{s,h0} = g_{s,h0} + \frac{1}{H''} \sum_{h=h_1}^{h_2} (gc_h - g_{s,h}) \quad (32)$$

$$W' = \max(p_4, W) \quad (33)$$

$$W = \sum_{s=1}^{N''} w_s \quad (34)$$

In Eqs. (31-34)  $N''$  stands for the number of partner series,  $h_1$  and  $h_2$  are thresholds of a time window around  $h0$ , and  $H''$  denotes the number of values used within that time window. Interpolated values of either the candidate or partner series are excluded in (32). All of  $h_1$ ,  $h_2$  and  $H''$  are function of both  $s$  and  $h0$  (see step 5.2), as well as the way of their determination varies according to the phase of the homogenization procedure. Weight  $w_s$  is a function of the spatial correlation ( $r$ ), the frequency of missing data around  $h0$  (as it may influence  $h_1$ ,  $h_2$  and  $H''$ ), and the status of observed data, i.e. they are within the homogenized period or not.

When it is not declared in other way,  $p_4 = 0.4$ . It follows from Eq. (33) that when  $W$  is very low,  $gc_{h0}$  approaches to 0, which is the climatic normal value.

Eqs. (31-34) show the case of monthly homogenization, but the same formulas are applied also for the gap filling in daily scale. Note that in the daily homogenization of any kind of variable the gap filling is always performed with data of daily resolution.

### B13 Gap filling for precipitation data

In this operation, data without logarithmic transformation are used, and the annual cycle of observed data is kept (35, 36).

$$\overline{xc_{[h_1, h_2]}} = \frac{1}{H''} \sum_{h=h_1}^{h_2} xc_h \quad (35)$$

$$\overline{xs_{[h_1, h_2]}} = \frac{1}{H''} \sum_{h=h_1}^{h_2} xs_{s,h} \quad (36)$$

Omitting  $[h_1, h_2]$  from the indexes of the long-term averages, and using  $W$  according to Eqs. (33-34), the solution is shown by (37). The sum of weights  $W$  below is defined by Eq. (34).

$$\begin{aligned} \text{if } W \geq p_4, \quad xc_{h0} &= \frac{1}{W} \sum_{s=1}^{N''} \frac{w_s x_{s,h0} \overline{xc}}{\overline{x_s}} \\ \text{if } W < p_4, \quad xc_{h0} &= \frac{1}{p_4} \left( \sum_{s=1}^{N''} \frac{w_s x_{s,h0} \overline{xc}}{\overline{x_s}} + (p_4 - W) \overline{x_m} \right) \end{aligned} \quad (37)$$

In Eq. (37),  $m$  stands for the calendar month to which  $h0$  belongs.

## IV PREPARATORY STEPS

From now the operation of ACMANTv4 is described in its logical and operational order. The operations are organised into 27 main steps of 5 sections. The main steps are often divided into sub-steps.

The sections describing the steps often start with basic information regarding the target variables, period of time series, time resolution, etc. However, such pieces of information are excluded when they are obvious from the section titles. Beyond this, target variables are not mentioned when the step is performed for all variables treated by ACMANTv4. For steps including relative time series construction or the use of relative time series, the working period (i.e. the period for which the operations are performed) is not mentioned, as it must be the homogenized period. Note that when a new time series is constructed from another one, the new series will always cover the whole series (which equals with the target period), even when the working period is shorter. Values outside the working period are usually copied to the new series, while missing data code is applied when the former option would not provide sensible values for the new series.

### 1 Setting parameters

Firstly some characteristics of the input dataset must be introduced, such as the number of time series, the target period of the homogenization, etc. Users may use the inbuilt physical thresholds or define thresholds which characterise better the climate represented by the data (optional). The rules of parameter setting are described in the Manual.

### 2 Basic control of input data

#### 2.1 Control of date order and input data format, and justification to the target period

Data out of the target period are removed. For series shorter than the target period, the missing sections are added with blocks of missing data codes. From this step, all the series have the same extent over the same period, which is the target period.

If the program finds a date order error or data formatting error, then it will stop with an error message.

#### 2.2 Control of physical outlier values

Section of time series: whole series (target period).

If the program finds a physical outlier, it will be considered missing data in the continuation.

### 3 Construction of networks

The dataset is divided into smaller networks when the number of time series is higher than 40, or when the spatial correlation ( $r^*$ ) is smaller than 0.4 for at least one pair of time series. Steps 3.1 – 3.4 are always done to check the spatial correlations (in step 3.5).

#### 3.1 Calculation of monthly values from daily values

Performed in daily homogenization. Section of time series: whole series.

Monthly values are calculated from daily values. The status of each monthly value is determined, they can be “observed” or “missing”.

#### 3.2 Determination of the treated period

#### 3.3 Transformation from monthly RR to monthly TR

Performed in RR homogenization. Section of time series: treated period.

Monthly RR is transformed to monthly TR (see B1).

#### 3.4 Calculation of deseasonalised monthly values

Performed for additive variables and for TR. Section of time series: treated period. Time resolution: monthly.

Climatic mean monthly values are subtracted from the observed values (see B2).

#### 3.5 Calculation of spatial correlation

Section of time series: treated period. Time resolution: monthly.

Spatial correlation ( $r^*$ ) is calculated for each pair of series of the dataset, according to B3.

#### 3.6 Network construction

Performed when  $N > 40$  or when  $r^* < 0.4$  for at least one pair of time series.

A distinct network is constructed for each series, in which they will be central series. The inclusion of highly correlated partner series and an even coverage of the central series with observed data of the partner series are favoured, while the increase of the number of time series in network over a limit is penalized. Data gaps or lapses in the

periods of observed data can reduce the number of the truly comparable observed data even in large size networks, therefore the concept of the number of effective partners ( $J$ ) is introduced, which indicates the number of synchronous observed values in partner series for any month of the central series.

This automatic networking has an important impact on the rest of the homogenization tasks, namely the homogenization results of the central series must be provided only, as each series of the dataset is central series in one network. By contrast, in 1-network homogenization, the homogenization results for all the series must be provided. In most steps of the homogenization procedure this difference does not appear, as ACMANT is based on the joint improvement of data quality and homogeneity in all time series within network. The few exceptions will be indicated.

Note that once the network construction is finished, data of different networks will never treated together within the homogenization procedure.

### 3.6.1 Selection of the most highly correlated partner series up to 30 series.

When the number of potential partner series with  $r^* \geq 0.4$  is higher than 30 for a given central series, the following steps (3.6.2 and 3.6.3) are performed recursively, while the networking for that central series terminates with this step in the reverse case.

### 3.6.2 Calculation of scores indicating the potential usefulness of further partner series

Scores are calculated for each of the series sufficiently correlating with the central series and not yet having been selected to be partner series.

i) Possible occurrences of having less than 10 partner series ( $J_h < 10$ ) is checked for every monthly value ( $h$ ). If such cases are found, the score  $S_1$  is calculated (38-39), otherwise  $S_1 = 0$ .

$$S_1(s) = 5 \sum_{h=1}^H r_s^{*4} (12 - J_h^*(s))^3 \quad (38)$$

$$J^* = \begin{cases} J & \text{if } J < 10 \\ 12 & \text{if } J \geq 10 \end{cases} \quad \text{for every } s \text{ and } h \quad (39)$$

ii) Ratios of cases  $J_h < 20$  are checked for all possible overlapping 10-year periods starting on January (referred to as decades). If decades with higher than 25% ratio of  $J_h < 20$  occur, score  $S_2$  is calculated (40-41), otherwise  $S_2 = 0$ . The cluster of dates belonging to at least one decade with >25% ratio of  $J_h < 20$  is denoted with  $Y$ .

$$S_2(s) = 5 \sum_{h \in Y} r_s^{*4} (20 - J_h^{**}(s))^2 \quad (40)$$

$$J^{**} = \begin{cases} J & \text{if } J < 20 \\ 20 & \text{if } J \geq 20 \end{cases} \quad \text{for every } s \text{ and } h \quad (41)$$



iii) Network size is scored with  $S_3$  (42).

$$S_3 = -(N' - 31)^2 \quad (42)$$

iv) Overall score ( $S$ ) is calculated for each  $s$  (43).

$$S(s) = S_1 + S_2 + S_3 \quad (43)$$

### 3.6.3 Selection of an additional partner series

The series with maximal score ( $S^* = \max\{S(s)\}$ ) is selected. If  $S^* > 0$ , the series of  $S^*$  will be a partner series, and the procedure continues from step 3.6.2. If  $S^* \leq 0$ , no series is selected at this step, and the networking for the given central series terminates here.

## 4 Preliminary operations

Steps 3.1 – 3.4 are repeated here with little difference. The reason of this repetition is that while the work was with the entire dataset in main step 3, from this step the work is always within a network.

### 4.1 Calculation of monthly values from daily values

The same as step 3.1.

### 4.2 Determination of possible excluded years

Section of time series: whole series.

### 4.3 Determination of the treated period

The same as step 3.2.

### 4.4 Transformation from monthly RR to monthly TR

The same as step 3.3.

### 4.5 For additive variables and for TR: Deseasonalisation

Performed for additive variables and for TR. Section of time series: whole series. Time resolution: monthly and daily.

The monthly climatic means are removed from the observed data according to B2 ( $\mathbf{X} \rightarrow \mathbf{G}$ ).

Note that in precipitation homogenization the  $\text{RR} \rightarrow \text{TR}$  transformation of monthly, annual and bi-seasonal values is repeated several times during the homogenization procedure. After such transformations the monthly TR values are always deseasonalised, while annual and bi-seasonal TR values never are.

## 5 Infilling data gaps

### 5.1 Calculation of spatial correlations

Section of time series: treated period. Time resolution: monthly.

Spatial correlation ( $r$ ) is calculated for each pair of monthly  $\mathbf{G}$  according to B3.

### 5.2 Gap filling

Section of time series: treated period. Time resolution: monthly or daily.

It is performed for each missing data of each series of the network, one-by-one. The series whose missing data are under the process of infilling is referred to as candidate series.

The basic formulas of gap filling are described in B12 – B13, but details as how many partner series are included or how wide time windows are used are not shown there. In the method presented here the use of relatively narrow time windows are preferred aiming to reduce the effect of possible inhomogeneities. However, the use of wider windows is allowed when there is not enough comparable data in narrower windows. For this reason, the window width around the missing data is specific both for the missing data position and for the pair of candidate series - partner series. Potential partner series are examined according to the decreasing order of  $r$ , but series without observed data for  $h_0$  ( $d_0$ ) are excluded.

Let the date of a missing data in the candidate series be  $h_0$  ( $d_0$  in daily homogenization). Firstly a relatively narrow symmetric time window is used around the year of  $h_0$  referred to as central year ( $y_0$ ). Pairs of observed data for series  $gc$  and a potential partner series  $s$  ( $xc$  and  $s$  for RR) are searched first in the central year, then in gradually increasing distance from  $y_0$ . The use of a given time window ( $p_5$ ) terminates when either the borders of the window are reached, or the required number of monthly value pairs or that of daily value pairs have been found. If the number of pairs of observed data is low within a given time window, then a wider  $p_5$  will be applied. When the required number of data pairs have been found, their required statistics are retained, and the procedure continues with another partner series until the correlations are sufficient ( $r = r_{gc,s} \geq 0.4$ ).

In the calculation of the interpolated value, the statistics of partner series determined by the use of narrower time windows are preferred in two ways: a) Potential partner series are ordered according to the used  $p_5$ , starting from the narrowest window. Partner series related to wider time windows are considered only when the number of partner series of narrower  $p_5$  is lower than threshold  $p_6$ , while the consideration of

further partner series terminates immediately when the number of partner series reaches threshold  $p_7$ . b) Weight ( $w_s$ ) of Eqs. (31) and (37) depends on the window width related to  $s$ .

Four time windows can be applied at this step. The related parameters and the weighting of partner series are shown in Table 1.

**Table 1.** Time windows and weights in the use of potential partner series  $s$  at step 5.2.  
 $p_5$  – width of time window around the central year;  $p_6$  – minimum number of partner series to accept a given time window;  $p_7$  – number of partner series at which the search of further partner series is finished;  $r$  – spatial correlation;  $w_s$  – weight of the observed value of series  $s$ .

$p_5$ (years)	Required number of data pairs		$p_6$	$p_7$	$w_s$
	Monthly	Daily			
7	60	1800	15	15	$r^2$
13	30	900	10	15	$0.9r^2$
25	30	900	5	15	$0.8r^2$
Unlimited	30	900	-	15	$0.5r^2$

Note 1: Numbers of monthly value pairs are counted also in gap filling of daily data, and the search of data pairs for a given  $s$  terminates when the number of either the monthly value pairs or daily value pairs reaches the prescribed threshold.

Note 2: If  $30 \leq H'' < 60$  monthly data pairs were found with  $p_5 = 7$  years,  $p_5 = 13$  years is not applied, as the number of data pairs exceeds the requirement for  $p_5 = 13$ .

### 5.3 Calculation of monthly, bi-seasonal and annual values

Section of time series: treated period.

When sums or arithmetical averages are calculated for larger time units than the source data resolution, observed values and interpolated values are treated equally.

#### 5.3.1 Calculation of monthly values

Performed in daily homogenization.

i) For additive variables: calculation of monthly means (arithmetical averages) from daily values.

ii) For precipitation: calculation of monthly total from daily values

#### 5.3.2 Calculation of annual and bi-seasonal values

i) In homogenization of additive variables: Calculation of annual means (arithmetical averages) from monthly values.

ii) In homogenization of additive variables with sinusoid seasonality: Calculation of annual values of summer – winter difference.

iii) In RR homogenization: calculation of annual totals from monthly values where precipitation is mostly rain throughout the year, and calculation of bi-seasonal totals where year is divided to a rainy season and a snowy season.

#### 5.4 Transformation from RR to TR

Monthly and annual or bi-seasonal RR values are transformed to TR values according to B1.

## V FIRST ITERATION

At this phase, the inhomogeneities of reference series have not reduced yet, hence the risk of making errors in the detection and adjustments is elevated. The aim is to detect relatively large breaks only, and to minimise the risk of unnecessary or too large adjustments. For this aim, strict significance thresholds are applied, and the probability of type one errors is kept low. In addition, the minimum adjustment terms of an ensemble homogenization is applied to reduce more the risk of overshooting.

### 6 Creation of relative time series for outlier filtering

Performed in monthly homogenization always, except for RR when user has switched off the outlier filtering. Performed in daily homogenization for additive variables. Not performed for daily RR. Time resolution: monthly.

Type of candidate series: **Gc**. Type of reference composites: **G**.

#### 6.1 Calculation of spatial correlations

Spatial correlations ( $r^*$ ) are calculated according to B3.

#### 6.2 Determination of sets of reference composites

Rules of B5 on the creation of sets of reference composites are applied. If  $N^* > 10$  for a given **F'** after the selection of reference composites according to B4, the reference composites are ordered according to  $r^*$ . Then the ordered reference composites are taken one-by-one and they are retained as far as each monthly value of the candidate series can be paired with at least 10 synchronous observed values of the retained reference composites. When this condition is fulfilled, the other reference composites are excluded.

#### 6.3 Weighting of reference composites

If  $N^* > 5$ , the weights of reference composites are calculated by ordinary kriging with some modifications (see step 6.3.1), while they are the squared spatial correlations ( $r^{*2}$ ) in the opposite case.

##### 6.3.1 Modified ordinary kriging

Firstly, the weights of reference composites ( $w_{s'}$ ) are calculated with the standard procedure of ordinary kriging, but these weights are sometimes modified and the final weights are denoted with  $w_s$ . Usually  $w_s = w_{s'}$  for all the reference composites, but modifications are applied when i)  $w_{s'} < 0$  for any  $s$ , or ii)  $w_{s'} > 0.4W$  for any  $s$ . The reasoning of (i) is that a negative weight does not have climatological interpretation within an area of spatially similar climate, while that of (ii) is that too large individual

weights are undesired, as any time series might include undetected errors or inhomogeneities.

i) if  $w_s' < 0$  then  $w_s = 0$ .

ii) if  $w_{s^*}' > 0.4W$  for series  $s^*$ , then  $w_s'$  values are increased by  $0.01w_{s^*}'$  for all the reference composites of  $w_s' < 0.4W$ . Then  $W$  is recalculated from the larger  $w_s'$  values. If relation  $w_{s^*}' > 0.4W$  is still true, the increase of the other weights will be repeated until  $w_s' > 0.4W$  is false for each  $s$ .

#### 6.4 Calculation of relative time series

After **F** has been created, the calculation of **T** is generally straightforward by B4, and also at this step, for additive variables. However, for monthly precipitation homogenization 3 kinds of relative time series are generated with 3 versions of the same candidate series. The relation between the data of the original candidate series and those of the other versions (indexed with  $a$  and  $b$ ) is shown by (44-45) for the RR data before transformation.

$$xc_{a,h} = \max(0, xc_h - 30) \quad \text{for every } h \quad (44)$$

$$xc_{b,h} = xc_h + 30 \quad \text{for every } h \quad (45)$$

Note that before the calculation of relative time series, all RR values are transformed to TR (as usual via B1), with which **Gc<sub>a</sub>** and **Gc<sub>b</sub>** are obtained. Finally, the relative time series will be denoted with **Ta** (**Tb**) for **Gc<sub>a</sub>** (**Gc<sub>b</sub>**).

Sets of relative time series for each candidate series are created according to B5. For monthly precipitation, such sets are created for each of **Ta** and **Tb**.

### 7 Filtering of spatial outliers

Time resolution: monthly.

The aim is to remove spatial outliers, but keep very low the risk of removing true extreme values. Therefore this routine is never performed for individual daily values. In RR homogenization the application of outlier filtering is more restricted than for the other climatic variables, as in RR data the frequent occurrence of dry days reduces the effective sample size influencing the signal-to-noise ratio.

Outlier filtered series (outlier filtered and deseasonalised series) will be referred to as **X<sup>+</sup>** (**G<sup>+</sup>**). For sections out of the homogenized period, or when this step is not applied,  $x^+ = x$  and  $g^+ = g$  for all the relevant dates.

Here, a few of the operations will be performed for all of the  $M$  relative time series of a candidate series. Therefore, best fitting relative time series (according to B6) will be denoted with **T<sub>i</sub>**.

## 7.1 Filtering of monthly spatial outliers

Performed for monthly homogenization except when user has switched off the outlier filtering. Not performed in daily homogenization.

7.1.1 The optimal  $\mathbf{T}_i$  for each year of the candidate series is selected according to B6.

### 7.1.2 Flagging likely outlier values

Performed for additive variables.

Standard deviation ( $\sigma$ ) is calculated for 3-month seasons of  $\mathbf{T}_i$ . Dates, for which (46) is true, are flagged as positions of possible outliers.

$$|t_{i,y,m} - \overline{\mathbf{T}_{i,m^*}}| > 5\sigma_{m^*} \quad (46)$$

### 7.1.3 Confirmation of an outlier value

Performed for additive variables.

A 19-month time window is selected around the flagged outlier  $t_{i,h}$ . The mean and standard deviation within this window are calculated excluding the flagged date, and the qualification as outlier is confirmed if (47) is true, while it is withdrawn in the opposite case.

$$|t_{i,h} - \overline{\mathbf{T}_{i,[h-9,h+9]}}| > 4\sigma_{[h-9,h+9]} \quad (47)$$

The exceedance of the starting or ending date of the homogenized period by a time window is not allowed. Therefore, when  $h$  is closer than 9 months to the first or last month of the homogenized period, the relevant half of the window will be truncated, while the other half remains 9-month wide.

### 7.1.4 Temporary adjustments for monthly outliers in relative time series

Performed for additive variables.

Confirmed outliers will be substituted with climatic averages in all of the relative time series (48).

If  $y$  and  $m$  define the date of an outlier, then

$$t'_{j,y,m} = t_{j,y,m} + \overline{\mathbf{T}_{i,m^*}} - t_{i,y,m} \text{ for every } j \in (1, 2, \dots, M) \quad (48)$$

The purpose of this adjustment is to provide the set of  $\mathbf{T}'$  without the detected monthly outliers for the next step (7.2, filtering of outlier periods).

### 7.1.5 Detection of outliers in precipitation series

$t_{i,y,m}$  is an outlier if any of relations (49-50) is true.

$$ta_{i,y,m} > 5\sigma_{m^*} \quad (49)$$

$$tb_{i,y,m} < -5\sigma_{m^*} \quad (50)$$

See the definition of **Ta** and **Tb** at step 6.4. The comparison of Eqs. (49-50) with Eq. (46) shows that the possibility of detecting precipitation outliers is restricted particularly for low precipitation totals.

### 7.1.6 Data quality is indicated as outlier for confirmed outliers.

## 7.2 Filtering of outlier periods

Performed for additive variables.

Concept: If the difference for a short section mean is large in comparison with the average value of adjacent sections of  $\mathbf{T}_i'$ , it is an indication of an outlier period. The significance is evaluated on normalized series ( $\mathbf{T}_i''$ ) to reduce the possible effects of seasonally varying variance. The detection of outlier periods is a step-by-step procedure, only the most significant outlier period is detected in a specific step. In certain later steps of the homogenization procedure, the endpoints of outlier periods longer than 4 months are treated as two independent breaks. By contrast, 1-4 month long outlier periods are always treated as outliers. In daily homogenization, the refinement of outlier periods (step 25) will examine these outliers in daily scale.

### 7.2.1 Calculation of normalized relative time series

The normalization is performed for all relative time series (51).

$$t_h'' = \frac{t_h' - \overline{\mathbf{T}_m'}}{\sigma_{m^*}} \quad (51)$$

### 7.2.2 Calculation of anomalies for sections of $\mathbf{T}''$

Performed for all relative time series.



Let  $h_1, h_2, h_3$  and  $h_4$  be four ordered points of series  $\mathbf{T}''$ . The size of the anomaly ( $q'$ ) of the central section  $[h_2+1, h_3]$  is the difference of its mean from the mean of the adjacent sections (52).

$$q' = \overline{\mathbf{T}''_{[h_2+1, h_3]}} - \frac{(h_2-h_1)\overline{\mathbf{T}''_{[h_1+1, h_2]}} + (h_4-h_3)\overline{\mathbf{T}''_{[h_3+1, h_4]}}}{h_2-h_1+h_4-h_3} \quad (52)$$

Let the length of the period is defined by (53).

$$L' = h_3 - h_2 \quad (53)$$

Note that  $[h_2+1, h_3]$  is not always the final position of the detected outlier period, that is why the symbols of the calculated statistics hold apostrophes.

Generally,  $q'$  is examined for all pairs of  $h_2$  and  $h_3$  for which  $0 < L' \leq 36$ , with three exceptions, a) in monthly homogenization  $1 < L' \leq 36$ , b) in mode switched off outlier filtering of monthly homogenization  $4 < L' \leq 36$ , c) note that at step 7.2.11 some value pairs will be excluded. While  $L'$  varies between 1 and 36 months, the length of the two adjacent sections ( $L^*$ ) is 24 months for each, at least when the distance from the endpoints of  $\mathbf{T}''$  allows that (see also steps 7.2.3 and 7.2.4).

### 7.2.3 Selection of optimal $\mathbf{T}_i''$ to each year of the candidate series

Usually the selection of  $\mathbf{T}_i''$  is performed according to B6. However, if the closeness of one endpoint of series  $\mathbf{T}_i''$  does not allow to edit 24-month sections for both  $[h_1+1, h_2]$  and  $[h_3+1, h_4]$ , then two cases are possible:

a) If period  $[h_2+1, h_3]$  can be examined with a relative time series whose endpoints are both at least 24-month distance from  $[h_2+1, h_3]$ , then a relative time series completing this condition will be used. If more than one relative time series complete the condition,  $\mathbf{T}_i''$  is selected from them according to B6. b) When period  $[h_2+1, h_3]$  is closer than 24 months to an endpoint of the homogenized period of the candidate series, the previous option does not work, and then the adjacent sections around  $[h_2+1, h_3]$  will be truncated, see next step.

### 7.2.4 Truncation of time window around the period examined

If the distance of  $h_2$  from the starting point of the homogenized period is shorter than 24 months (12 months),  $L^*(h_1, h_2)$  equals 12 months (zero), and the changes follow the same logic when  $h_3$  is closer to the endpoint of the homogenized period than 24 months. If  $L^* = 0$  for one of the adjacent sections, then the other adjacent section will be extended to 36 months.

### 7.2.5 Significance of anomalies

#### i) Significance of anomalies without consideration of seasonal cycle

The significance ( $\alpha$ ) is a function of the anomaly and the length of the period (54).

$$\alpha = (L')^{0.8}(q')^2 \quad (54)$$

Note that the theoretical solution for temporally independent variables would be similar to Eq. (54) but without exponent. With the exponent  $< 1$ , the possible autocorrelation of relative time series is taken into consideration.

#### ii) Significance of anomaly with consideration of seasonal cycle

The significance increases with duration (54), however, in case of sinusoid model, it is taken into consideration that the same season deviations from the mean of the adjacent years might indicate seasonal variation instead of altered value of the means. Therefore,  $L^\#$  ( $L^\# \leq L'$ ) is applied in Eq. (54) instead of  $L'$ , and  $L^\#$  is calculated according to (55):

$$L^\# = \max\{1, L' - \frac{0.75}{3.5} \sum_{h=h_2+1}^{h_3} \mu_m\} \quad (55)$$

In univariate homogenization  $L^\# = L'$ .

### 7.2.6 Pre-selection and further examination of periods with significant anomalies

Periods with  $\alpha \geq 25$  are pre-selected. A supplementary examination is performed when  $L^* > 0$  both for  $[h_1+1, h_2]$  and  $[h_3+1, h_4]$ , i.e. Eqs. (52-55) are repeated using only one adjacent section of  $[h_1+1, h_2]$  and  $[h_3+1, h_4]$ .  $q$  and  $\alpha$  will be indexed with A (B) when  $q'$  is calculated with the left (right) hand side adjacent section. A pre-selected period of  $\alpha \geq 25$  will be retained, if all of the following relations are true (56).

$$\alpha_A \geq 25; \quad \alpha_B \geq 25; \quad \text{sign}(q_A) = \text{sign}(q_B) \quad (56)$$

When  $L^* = 0$  for any of  $[h_1+1, h_2]$  and  $[h_3+1, h_4]$ , the pre-selected period is always retained.

### 7.2.7 Final selection of the period with the most significant anomaly

When all the possible periods have been examined, the one with the maximal  $\alpha$  will be selected from those which have been retained in step 7.2.6. If no period has been retained in step 7.2.6, the procedure of outlier period filtering terminates for a given candidate series.

### 7.2.8 Refinement of the starting end ending months of outlier period

Performed for sinusoid cycle of inhomogeneities.

Data within and around the outlier period selected at 7.2.7 is examined further. The optimal modified step function is fitted to the  $t_i''$  values of  $[h_1+1, h_4]$ , according to B9.  $h_2$  and  $h_3$  are allowed to change to  $h_2^*$  and  $h_3^*$  according to (57-58).

$$h_2 - 14 < h_2^* \leq h_2 \quad (57)$$

$$h_3 \leq h_3^* < h_3 + 14 \quad (58)$$

As the timings of the two endpoints of the outlier period are searched, usually the number of breakpoints  $K = 2$ , but if  $L^* = 0$  either for  $[h_1+1, h_2]$  or for  $[h_3+1, h_4]$ , then  $K = 1$ . When  $h_3^* - h_2^* < 10$ , the new estimations are discarded and  $[h_2^*+1, h_3^*] \equiv [h_2+1, h_3]$ .

If  $[h_2^*+1, h_3^*]$  differs from  $[h_2+1, h_3]$ , then the mean bias and length are recalculated (Eqs. 52-53) with the new parameters, and they are denoted with  $q$  and  $L$ , respectively. If  $[h_2^*+1, h_3^*] \equiv [h_2+1, h_3]$ , then  $q = q'$  and  $L = L'$ .

### 7.2.9 Temporal adjustments

In step 7.2.8.  $h_2^*$ ,  $h_3^*$ ,  $L$  and  $q$  were determined for sinusoid cycle of inhomogeneities. If the model seasonality is not sinusoid, then  $[h_2^*+1, h_3^*] \equiv [h_2+1, h_3]$ ,  $q = q'$  and  $L = L'$ .

$q$  is subtracted from all the values of  $\mathbf{T}_{[h_2+1, h_3]}''$  in all relative time series belonging to the given candidate series. This is necessary for the search of further outlier periods, but these adjusted values are not transmitted to other subroutines of the homogenization procedure.

### 7.2.10 Indication of data quality within detected outlier periods

When  $L \leq 4$ , the data are indicated as outliers. When  $4 < L \leq 28$ , the data are indicated as part of a short-term inhomogeneity. When  $L > 28$ , data quality remains without indication. The only purpose of the search and temporary correction of inhomogeneities of 29-36 months in the procedure of outlier filtering is to provide a more accurate detection for the shorter outlier periods.

The observed data in daily series (when they exist) receive the same quality indication as the indication of monthly data. (Note that this rule is valid for outlier filtered values, but not true for missing values.)

### 7.2.11 Exclusion of flip-flop

Rarely it happens that after the temporal adjustment of a section, its adjacent section (or a part thereof) seems to be outlier, and after the adjustment of the latter, the former

becomes outlier again. This could lead to an infinite cycle, and to avoid it,  $L'$  of further outlier periods ( $L^+$ ) is maximised by  $0.8L'$  for any outlier period could be detected later around  $h_2$  (59).

$$L^+(h_2^+) < 0.8L'(h_2) \text{ for any } h_2^+ \in [h_2 - 0.8L, h_2 + 0.8L] \quad (59)$$

## 8 Infilling data gaps

This main step is performed when any of steps 7.1 and 7.2 is performed. Section of time series: treated period.

Detected outliers or values belonging to detected short-term inhomogeneities are treated in the same way as the missing data, and they are substituted with interpolated values.

8.1 Calculation of spatial correlations ( $r$ ) for each pair of monthly  $\mathbf{G}^+$  series.

8.2 Gap filling

The same as step 5.2, except for some tiny differences detailed in (i) and (ii).

i)  $p_6$  and  $p_7$  have different parameterization, while the other parameters are unchanged (see Table 2).

**Table 2.** Time windows and weights in the use of potential partner series  $s$  at step 8.2 (Symbols are explained at Table 1).

$p_5$ (years)	Required number of data pairs		$p_6$	$p_7$	$w_s$
	Monthly	Daily			
7	60	1800	7	10	$r^2$
13	30	900	5	10	$0.9r^2$
25	30	900	2	10	$0.8r^2$
Unlimited	30	900	-	10	$0.5r^2$

ii) Data considered in Eqs. (31), (35) and (36) are restricted to the homogenized period when  $p_5 = 7$  years, while no such restriction is applied when the window is wider.

### 8.3 Calculation of monthly, bi-seasonal and annual values

In daily homogenization, monthly values are calculated both for  $G$  and  $G^+$  (for additive variables) or for  $X$  and  $X^+$  (for RR). Annual and bi-seasonal values are calculated only for  $G^+$  or  $X^+$ . The calculations are made in the same way as in step 5.3.

### 8.4 Transformation from RR to TR

The same as step 5.4.

## 9 Creation of relative time series for break detection

Time resolution: annual or bi-seasonal. Type of candidate series:  $Gc^+$ . Type of reference composites:  $G^+$ .

### 9.1 Calculation of spatial correlations ( $r^*$ )

This step is performed with monthly data, as usual.

### 9.2 Determination of sets of reference composites

Rules of B5 are applied for the annual variables detailed at step 5.3.2. All the adequate reference composites are included.

### 9.3 Weighting of reference composites: they are equally weighted at this step.

### E9.4 Truncation of sets of reference composites

From this step, an ensemble cycle starts for the homogenization of annual variables. Steps 9.4 – 11.2 are performed for each candidate series of the network, and always with the exclusion of 1 partner series, which can be done in  $N'-1$  ways, hence the number of the repetition of the ensemble cycle is  $N'(N'-1)$ , and the number of the ensemble members for each candidate series is  $N'-1$ .

With 9.2 - 9.3, the sets of reference series have been constructed, however, as one partner series is excluded from this step, the sets of reference composites are truncated. With this truncation the number of reference composites may fall below 3, but it is allowed within the ensemble cycle.

E9.5 Calculation of relative time series according to B4.

## 10 Break detection

Time resolution: annual.

### E10.1 Selection of relative time series

The principal rule is to select the series with the highest  $\beta$ , as it is described at B6. However, to reducing edge effects, a 10-year overlapping is applied. The operation of this overlapping is shown here with an example: Let suppose that the break detection for section  $(y_1, y_2)$  of the candidate series has been performed with  $\mathbf{T}_A$ , and sections before  $y_1$  and after  $y_2$  can be examined with other relative time series. In the continuation,  $\mathbf{T}_B$  is involved for the break detection in section before  $y_1+10$  or in section after  $y_2-10$  or in both. However, if a break has been detected with  $\mathbf{T}_A$  before  $y_1+10$  or after  $y_2-10$ , the length of overlapping is reduced until the date of the detected break.

### E10.2 Break detection with step function fitting

The minimal distance between two steps is 3 years by definition.

i) Break detection for additive variables of non-sinusoid seasonality

Univariate detection (B7) is applied with  $p_2 = 3.92$ .

ii) Break detection for additive variables of sinusoid seasonality

Bivariate detection (B8) is applied where the two variables are the annual mean and the summer – winter difference, and the parameters are  $p_2 = 2.8$  and  $p_3 = 0.2$ .

iii) Break detection for precipitation without seasonal difference of the form of precipitation

The same as case i).

iv) Break detection for precipitation with rainy season and snowy season

Bivariate detection (B8) is applied where the two variables are the TR of rainy season and the TR of snowy season, and  $p_2 = 2.8$ .  $p_3$  depends on the length of the snowy season ( $L^{(S)}$ , in months) according to (60).

$$p_3 = \left( \frac{L^{(S)}}{12 - L^{(S)}} \right)^2 \quad (60)$$

### E10.3 Control with t-test

When bivariate detection is applied, this step is performed for each of the two variables.

Common t-test is applied for checking the significance of breaks detected in step E10.2. For break  $k$ , the difference of section means in  $\mathbf{T}_j$  between the periods  $[y_{k-1}, y_k]$  and  $[y_k, y_{k+1}]$  is tested.  $\mathbf{T}_j$  stands for the relative time series in which the break was detected. The standard deviation ( $\sigma$ ) of the t-test is the  $\sigma$  for the entire period of  $\mathbf{T}_j$  in order to reduce estimation errors which could come from statistics of short sections.

Critical values applied facilitate 0.01 probability of first type error in general and 0.007 probability of first type error for breaks of summer – winter difference ( $\tilde{\mathbf{T}}_j$ ). Non-significant breaks are removed from the break list.

### E10.4 Limitation of the number of synchronous breaks

When bivariate detection is applied, this step is performed for each of the two variables.

A high ratio of coincidental detected breaks gives instability to the homogenization results, and in the extreme case of synchronous breaks in all series of the network, the equation system of ANOVA correction model is indefinite. To limit this instability, the number of synchronous breaks ( $K^\#$ ) must be smaller than the half of the number of time series ( $N^\#$ ) whose homogenized periods include the date of synchronous breaks (61).

$$K_y^\# < 0.5N_y^\# \quad \text{for every } y \in (1, 2, \dots, n') \quad (61)$$

When condition (61) is violated, the breaks of the lowest significance ( $\alpha$ ) are omitted one-by-one until relation (61) becomes true. The significance is calculated by (62-63).

$$\alpha_{s,k} = \frac{(y_k - y_{k-1})(y_{k+1} - y_k)Q_{s,k}}{y_{k+1} - y_{k-1}} \quad (62)$$

$$Q_{s,k} = (\delta \bar{\mathbf{T}})^2 = (\overline{\mathbf{T}_{s[y_k+1, y_{k+1}]}} - \overline{\mathbf{T}_{s[y_{k-1}+1, y_k]}})^2 \quad (63)$$

## 11 Adjustments for inhomogeneities

### E11.1 Application of ANOVA correction model

Section of time series: homogenized period. Type of input series:  $\mathbf{G}^+$ . Time resolution: annual.

The ANOVA correction model (B10) is applied for the variable(s) examined in main step 10. The result will be a vector of annual adjustment terms ( $\mathbf{Z}^*$ ) (or vectors  $\mathbf{Z}^*$  and  $\mathbf{Z}^{**}$  in bivariate cases) for each time series of the network.

Note that B10 or B11 can be applied for variables like summer – winter difference in the same way as for annual, monthly or daily means.

## E11.2 Calculation of adjustment terms backwards from the beginning of the homogenized period

Section of time series: From the starting of the treated period until the starting of the homogenized period. Time resolution: annual.

This step is performed when the homogenized period starts later than the treated period, and the ANOVA correction model results in non-zero adjustment term for the first year of the homogenized period ( $x_{y_0+1}$ ). The idea behind the backwards adjustments is that a bias for inhomogeneities in the homogenised period may indicate a non-zero bias of the earlier sections of the time series, although with less certainty and accuracy. Biases are intended to be reduced here in the entire treated period, and in certain later steps in the entire time series.

Term “long-term adjustment term” ( $z_L$ ) is introduced here as the minimum of the adjustment term for  $x_{y_0+1}(z_A)$  and of the average adjustment term for the first 30 years of the homogenized period ( $z_B$ ), (64).

$$z_L = \begin{cases} \min(|z_A|, |z_B|) & \text{if } \text{sign}(z_A) = \text{sign}(z_B) \\ 0 & \text{if } \text{sign}(z_A) \neq \text{sign}(z_B) \end{cases} \quad (64)$$

The adjustment term ( $z^*$ ) for the section before  $x_{y_0}$  is identical with  $z_L$ , with the exception that when  $z_L \neq z_A$ ,  $z^*$  gradually changes from  $z_A$  to  $z_L$  over a 3-year period, going backwards from  $y_0$  to  $y_0-2$ . Note when the homogenized period is shorter than 30 years,  $z_B = 0$ .

## 11.3 Adjustment terms derived from ensemble results

Section of time series: treated period.

The ensemble cycle is terminated with the previous step, and the procedure follows with the evaluation of the ensemble adjustment terms.

### 11.3.1 Adjustment terms for the examined annual variables

Time resolution: annual.



For each year  $y$  of the treated period,  $z'_y$  will be equal with the minimum absolute value of the ensemble members for  $z_y^*$  (65).

$$z'_y = \begin{cases} \text{sign}(z_{1,y}^*) \min(|z_{1,y}^*|, |z_{2,y}^*| \dots |z_{N',y}^*|) & \text{if } \text{sign}(z_{a,y}^*) = \text{sign}(z_{b,y}^*) \\ & \text{for every } a, b \in [1, 2 \dots N']; \\ 0 & \text{if } \text{sign}(z_{a,y}^*) \neq \text{sign}(z_{b,y}^*) \text{ for any pair of } a, b \end{cases} \quad (65)$$

In this way, vector  $\mathbf{Z}'$  (in case of bivariate procedure, vectors  $\mathbf{Z}'$  and  $\mathbf{Z}''$ ) is constructed from  $\mathbf{Z}^*$  (from  $\mathbf{Z}^*$  and  $\mathbf{Z}^{**}$ ).

### 11.3.2 Adjustment terms of monthly and daily resolution

Time resolution: monthly and daily.

In univariate procedures the final adjustment terms for application ( $z$ ) are identical with the final adjustment terms for the examined annual variables ( $z'$ ) for any month and any day of a given year, but in bivariate procedures  $z$  depends on both  $z'$  and  $z''$ .

i) Additive variables with sinusoid cycle, monthly adjustment terms (66)

$$z_{y,m} = z'_y + 0.55 \sin\left(\frac{2\pi(m-2.7)}{12}\right) z''_y \quad (66)$$

Note: The constant in the numerator differs from that in Eq. (26), since for timings of annual (monthly) preciseness the timing is the last day of the year (month) by definition, while for adjustments the middle of the month represents best a month.

ii) Additive variables with sinusoid cycle, daily adjustment terms

Monthly adjustment terms are considered accurate for the middle day of the month.

Daily adjustment terms are determined with linear interpolation between the adjacent mid-monthly values.

iii) Precipitation with rainy season and snowy season (67-68)

$$z_{y,m} = \begin{cases} z'_y & \text{if } m \in \text{rainy season} \\ z''_y & \text{if } m \in \text{snowy season} \end{cases} \quad (67)$$

$$z_{y,m,d} \equiv z_{y,m} \text{ for every } d \quad (68)$$

## 11.4 Execution of adjustments

Section of time series: treated period. Time resolution: monthly and daily.

Both of  $\mathbf{G}$  and  $\mathbf{G}^+$  ( $\mathbf{X}$  and  $\mathbf{X}^+$ ) are adjusted for additive variables (for RR), and the adjusted series will be referred to as  $\mathbf{G}^*$  and  $\mathbf{G}^{+*}$  ( $\mathbf{X}^*$  and  $\mathbf{X}^{+*}$ ), respectively.

### i) Additive variables

In monthly homogenization, monthly values are adjusted. In daily homogenization, both of daily and monthly values are adjusted. Eqs. (69-70) are valid either for daily or monthly series.

$$\mathbf{G}^{+*} = \mathbf{G}^+ + \mathbf{Z} \quad (69)$$

$$\mathbf{G}^* = \mathbf{G} + \mathbf{Z} \quad (70)$$

### ii) Precipitation

In monthly homogenization, monthly RR series are adjusted, while in daily homogenization only daily RR series are adjusted. Eqs. (71-72) show the execution of daily adjustments.

$$x_{y,m,d}^{+*} = x_{y,m,d}^+ e^{z_{y,m,d}} \text{ for every } y,m,d \quad (71)$$

$$x_{y,m,d}^* = x_{y,m,d} e^{z_{y,m,d}} \text{ for every } y,m,d \quad (72)$$

## 11.5 Aggregation of daily values to create monthly values

Performed in daily RR homogenization.

Daily values are aggregated to create monthly values both for  $\mathbf{X}^*$  and  $\mathbf{X}^{+*}$ .

## 11.6 Transformation from RR to TR

In monthly RR homogenization, monthly RR values are transformed to TR according to B1.  $\mathbf{X}^* \rightarrow \mathbf{G}^*$ ,  $\mathbf{X}^{+*} \rightarrow \mathbf{G}^{+*}$ . Series of  $\mathbf{Ga}^*$  and  $\mathbf{Gb}^*$  are created according to step 6.4.

## VI SECOND ITERATION

The purpose of this homogenization cycle is twofold: a) improve accuracy with the use of pre-homogenized reference series, and b) manifest the range of instability of the homogenization results, i.e. when small changes in the relative time series construction or parameterization yield different homogenization results. For the latter purpose, the significance thresholds are light, and the probability of type one error is relatively high. The instability range appears as the difference between the minimum and average adjustment terms of an ensemble homogenization.

### 12 Creation of relative time series for outlier filtering

Time resolution: monthly. Candidate series:  $\mathbf{Gc}^*$  (for precipitation  $\mathbf{Gc_a}^*$  and  $\mathbf{Gc_b}^*$ ).  
Reference composites:  $\mathbf{G}^{+*}$ .

12.1 Calculation of spatial correlations ( $r^*$ )

12.2, 12.3, 12.4: The same as steps 6.2, 6.3, and 6.4.

### 13 Outlier filtering

It is the same as main step 7. Note that the data quality flags of the previous round of outlier filtering are cancelled before the execution of this main step.

### 14 Infilling data gaps

Section of time series: treated period.

14.1 Calculation of spatial correlations ( $r$ ) for each pair of monthly  $\mathbf{G}^{+*}$  series

14.2 Gap filling

It is performed similarly to the gap filling of step 5.2, with the following differences:

i) In the collection of data pairs used in Eq. (32) or Eq. (37), the values of both the candidate series and partner series must fall within the homogenized period. Note, however, that the period for infilling data gaps is still the treated period.

ii) Parameters  $p_5$ ,  $p_6$ ,  $p_7$  and  $w_s$  are changed (Table 3).

iii) The concept “type 1 series” is introduced for series whose homogenized period include  $y_0$  ( $y_0$  is used with the same meaning as in step 5.2) and “type 2 series” for the other series. When partner series  $s$  is a type 2 series, the weights are halved (Table S3).

iv) When the candidate series is type 1 series and series  $s$  is type 2 series, series  $s$  will be considered only if the number of partner series of type 1 would be less than 3, or the sum of the weights of the partner series of type 1 would be less than 0.4. When any of the latter two relations is true, type 1 partner series are retained, and time windows parameterized by  $p_5$ ,  $p_6$ ,  $p_7$  are examined again for finding possible partner series of type 2.

**Table 3.** Time windows and weights for potential partner series  $s$  at step 14.2.  
 $w_s(1)$  – the date of the missing data of the candidate series ( $h_0$  or  $d_0$ ) is within the homogenized period in the partner series;  $w_s(2)$  – the date of the missing data of the candidate series is out of the homogenized period in the partner series; (other symbols are explained at Table 1).

$p_5$ (years)	Required number of data pairs		$p_6$	$p_7$	$w_s(1)$	$w_s(2)$
	Monthly	Daily				
7	60	1800	7	10	$r^2$	$0.5r^2$
20	30	900	7	10	$0.92r^2$	$0.46r^2$
40	30	900	5	10	$0.85r^2$	$0.425r^2$
Unlimited	30	900	(3)	10	$0.7r^2$	$0.35r^2$

#### 14.3 Calculation of monthly, bi-seasonal and annual values within the treated period

In daily homogenization, monthly values are calculated for each of  $\mathbf{G}$ ,  $\mathbf{G}^+$  and  $\mathbf{G}^{+*}$  (for additive variables) or for each of  $\mathbf{X}$ ,  $\mathbf{X}^+$  and  $\mathbf{X}^{+*}$  (for RR). Annual and bi-seasonal values are calculated for  $\mathbf{G}^+$  and  $\mathbf{G}^{+*}$  or for  $\mathbf{X}^+$  and  $\mathbf{X}^{+*}$ . The calculations are made in the same way as in step 5.3.

#### 14.4 Transformation from RR to TR

The same as step 5.4.

### 15 Creation of relative time series for break detection

Time resolution is annual or bi-seasonal. Type of candidate series:  $\mathbf{Gc}^+$ . Type of reference composites:  $\mathbf{G}^{+*}$ .

15.1 Calculation of spatial correlations ( $r^*$ ) using  $\mathbf{G}^{+*}$  type series. This step is performed with monthly data, as usual.

15.2 Determination of sets of reference composites

Same as step 9.2.

E15.3 Weighting of reference composites

With this step, the same kind ensemble homogenization cycle starts as between steps E9.4 and E11.2, with the exclusion of one reference composite specific for the ensemble member. One difference here from the previous ensemble cycle is that the weighting of reference composites depends on the spatial correlations, therefore the weighting is part of the ensemble cycle.

When  $N^* > 6$ , the weights are calculated by ordinary kriging with the modifications described at step 6.3. When  $N^* = 3$  the weights are equal, while for  $3 < N^* < 6$ , the weights are the squared spatial correlations ( $r^{*2}$ ). Note, that for the exclusion of one reference composite the number of reference composites used is  $N^* - 1$ .

E15.4 Calculation of relative time series according to B4

## 16 Break detection

Time resolution: annual.

E16.1 Selection of relative time series

The same as step 10.1.

E16.2 Break detection with step function fitting

The same as step 10.2, except that  $p_2 = 1.4$  ( $p_2 = 1.0$ ) in univariate (bivariate) detection.

E16.3 Control with t-test

The same as step 10.3, except that the significance thresholds are lighter. The applied critical values of t-test allow first type error probability of  $\sim 0.3$  for  $\tilde{\mathbf{T}}_j$  and  $\sim 0.4$  for the other variables. Note that true first type error frequency is lower than these probabilities, as breaks are accepted only if both step function fitting and t-test confirm them.

#### E16.4 Limitation of the number of synchronous breaks

The same as step 10.4.

### 17 Adjustments for inhomogeneities

E17.1 – E17.2: The same as steps 11.1 – 11.2.

#### 17.3 Adjustment terms derived from ensemble results

The ensemble cycle has been terminated with the previous step, and the procedure follows with the evaluation of the ensemble members.

##### 17.3.1 Adjustment terms derived directly from ensemble results

Time resolution: annual.

Two kinds of adjustment terms constructed here: one is the same as determined by Eq. (65) ( $\mathbf{Z}'$  and  $\mathbf{Z}''$ ), while the other is the arithmetical average of the ensemble results, denoted with  $\mathbf{Z}^+$  in general and with  $\mathbf{Z}^{++}$  for the second variable in bivariate homogenization.

##### 17.3.2 Annual adjustment terms for 9 scenarios

Time resolution: annual.

Nine scenarios are created by the linear combination of  $\mathbf{Z}'$  and  $\mathbf{Z}^+$ . The serial number of ensemble member is denoted by upper index in brackets (73).

$$\mathbf{Z}^{(i)'} = p^{(i)}\mathbf{Z}' + (1 - p^{(i)})\mathbf{Z}^+, \quad i \in (1, 2, \dots, 9) \quad (73)$$

$p^{(i)}$  is of a Gaussian distribution with 0.5 expected value. The standard deviation is based on experiments of efficiency tests.  $p^{(1)} = -3.0$ ,  $p^{(2)} = -1.8$ ,  $p^{(3)} = -0.94$ ,  $p^{(4)} = -0.19$ ,  $p^{(5)} = 0.5$ ,  $p^{(6)} = 1.19$ ,  $p^{(7)} = 1.94$ ,  $p^{(8)} = 2.8$ ,  $p^{(9)} = 4.0$ . The scenario with  $p^{(5)}$  is referred to as most probable scenario.

In case of bivariate homogenization, the same relation is valid for  $\mathbf{Z}^{(i)''}$ ,  $\mathbf{Z}''$  and  $\mathbf{Z}^{++}$  as which is defined by Eq. (73) for the first variable. In the continuation, both  $\mathbf{Z}'$  and  $\mathbf{Z}^{(i)'}$  (and in bivariate homogenization also  $\mathbf{Z}''$  and  $\mathbf{Z}^{(i)''}$ ) will be in use for data adjustments. For distinguishing  $\mathbf{Z}'$  and  $\mathbf{Z}''$  from the ensemble adjustment terms, they will be referred to as standard adjustment terms.

## 17.4 Standard monthly adjustments

$\mathbf{Z}'$  and  $\mathbf{Z}''$  are used for standard monthly adjustments.

### 17.4.1 Calculation of standard monthly adjustment terms

In case of univariate homogenization,  $z_{y,m} = z'_y$  for any  $y$  and  $m$ . In bivariate homogenization Eqs. (66) and (67) are applied to determine the values of  $\mathbf{Z}$ .

### 17.4.2 Application of standard monthly adjustments

For additive variables,  $\mathbf{G}^{+*}$  and  $\mathbf{G}^*$  are determined by Eqs. (69) and (70), respectively. In RR homogenization only the outlier filtered series ( $\mathbf{X}^+$ ) are adjusted. Values of  $\mathbf{X}^{+*}$  are determined by Eq. (71).

### 17.4.3 Transformation from RR to TR

Monthly values of  $\mathbf{X}^{+*}$  are transformed to  $\mathbf{G}^{+*}$  according to B1.

## 17.5 Ensemble adjustments

$\mathbf{Z}^{(i) '}$  and  $\mathbf{Z}^{(i) ''}$  are used for determining monthly and daily adjustment terms.

### 17.5.1 Calculation of the ensemble adjustment terms

In univariate homogenization,  $z_{y,m}^{(i)} = z_y^{(i) '}$  for any  $i$ ,  $y$  and  $m$ , which means that the adjustment terms are constant within a calendar year. The uniformity of adjustment terms within a calendar year is valid also for daily adjustment terms.

In bivariate homogenization the values of  $\mathbf{Z}^{(i)}$  are constructed from  $\mathbf{Z}^{(i) '}$  and  $\mathbf{Z}^{(i) ''}$  in the same way as  $\mathbf{Z}$  is constructed from  $\mathbf{Z}'$  and  $\mathbf{Z}''$  in step 11.3.2.

### 17.5.2 Application of ensemble adjustments

For additive variables,  $\mathbf{G}^{(i)+*}$  is constructed from  $\mathbf{G}^+$  and  $\mathbf{Z}^{(i)}$  for each scenario ( $i$ ) by Eq. (69). In RR homogenization, the scenarios of  $\mathbf{X}^{(i)+*}$  are constructed from  $\mathbf{X}^+$  and  $\mathbf{Z}^{(i)}$  by Eq. (71).

## VII THIRD ITERATION

Based on the instability range calculated in the previous homogenization cycle, 9 scenarios of pre-homogenization results are taken, and the homogenization is done for each of these 9 scenarios. The homogenization result of this cycle will be the mean of the 9-member ensemble homogenization.

In this cycle some new steps appear, which were not present in the previous two homogenization cycles (e.g. downscaling break detection results to monthly and daily time scales), since they are needed more to the final accuracy of the results than to the success of the iteration.

### 18 Creation of relative time series for outlier filtering

This main step is performed for additive variables. Time resolution: monthly. Candidate series:  $Gc^*$ . Reference composites:  $G^{+*}$ .

This main step is performed in the same way as main step 6.

### 19 Outlier filtering

For additive variables this main step is the same as main step 7. Note that the data quality indications of the previous outlier filtering are cancelled before the execution of this main step.

In RR homogenization this main step is not performed, and the data quality indications of the previous outlier filtering (main step 13) are preserved.

### 20 Infilling data gaps and preliminary calculations for ensemble homogenization

With the infilling of data gaps an ensemble homogenization cycle starts. Before that, some preliminary calculations are performed. The spatial covariance and correlation matrix are determined by the use of standard adjustment terms, hence these characteristics, as well as the homogenized period remain the same in the ensemble homogenization.

#### 20.1 Calculation of spatial correlation

Both kinds of spatial correlations ( $r$  and  $r^*$ ) are calculated, for each pair of time series.



## 20.2 Determination of sets of reference composites

Two sets are calculated. For both sets the rules of B5 are applied, but in the second set  $N^*$  is limited in the same way as in step 6.2. The first set will be used in the operations with annual or bi-seasonal data, while the set of limited number of reference composites will be used in the operations with monthly data.

## 20.3 Weighting of reference composites

For both sets of reference composites: the same as step 6.3.

## E20.4 Gap filling

Section of time series: homogenized period.

In this step an ensemble homogenization cycle starts. The cycle is executed for the 9 scenarios generated in step 17.5.

This step is performed similarly to the gap filling in step 14.2. The concepts of series type 1 and type 2 are used in the same way as in the gap filling as in step 14.2 with some small differences, which are as follows.

- i) Gap filling is performed only for the homogenized period of time series.
- ii) Parameters  $p_5$ ,  $p_6$ ,  $p_7$  and  $w_s$  are changed relative to step 14.2 (Table 4).

**Table 4.** Time windows and weights for potential partner series  $s$  at step 20.4.  
 $w_s(1)$  – the date of the missing data of the candidate series ( $h_0$  or  $d_0$ ) is within the homogenized period in the partner series;  $w_s(2)$  – the date of the missing data of the candidate series is out of the homogenized period in the partner series; (other symbols are explained at Table 1).

$p_5$ (years)	Required number of data pairs		$p_6$	$p_7$	$w_s(1)$	$w_s(2)$
	Monthly	Daily				
25	100	3000	7	10	$r^2$	$0.5r^2$
51	30	900	5	10	$0.9r^2$	$0.45r^2$
Unlimited	30	900	(3)	10	$0.7r^2$	$0.35r^2$

## E20.5 Calculation of monthly, bi-seasonal and annual values within the homogenized period

The same as step 14.3.

## E20.6 Transformation from RR to TR

The same as step 5.4.

## E20.7 Time series without adjustments for short-term inhomogeneities

In some operations outlier periods of  $L \geq 5$  are considered short term inhomogeneities delimited by two breaks.  $\mathbf{G}^\#$  series are constructed here, in which the values of outlier periods of  $L < 5$  are substituted with interpolated values, but those of the longer outlier periods, as well as the inhomogeneities are left unchanged. Note that in RR homogenization filtering of outlier periods is not included, hence for TR data  $\mathbf{G}^\# \equiv \mathbf{G}^+$ .

## 21 Creation of relative time series

### E21.1 Relative time series for annual or bi-seasonal variables

Type of candidate series:  $\mathbf{Gc}^+$ . Type of reference composites  $\mathbf{G}^{+*}$ .

Rules of B4 are applied. The number of reference composites is unlimited (see also step 20.2).

### E21.2 Relative time series for monthly variables

Type of candidate series:  $\mathbf{Gc}^+$  and  $\mathbf{Gc}$ . Type of reference composites  $\mathbf{G}^{+*}$ .

Rules of B4 are applied. The number of reference composites is determined according to step 6.2. Relative time series for  $\mathbf{Gc}^+$  are denoted with  $\mathbf{T}$ , while those for the raw, no outlier filtered series ( $\mathbf{Gc}$ ) are denoted with  $\mathbf{T}^*$ .

## 22 Break detection

Time resolution: annual, monthly and daily.

### E22.1 Selection of relative time series

The same as step 10.1.

## E22.2 Break detection with step function fitting

Time resolution: annual.

The same as step 10.2, except that  $p_2 = 2.8$  ( $p_2 = 2.0$ ) in univariate (bivariate) detection.

## E22.3 Audit of outliers connected with breaks

Time resolution: monthly. Not performed in daily precipitation homogenization.

When a detected break is connected to adjacent detected monthly outlier values, it is controlled if these months are a part of the detected long-term inhomogeneity, or they are true outliers. This control is performed for each detected break.

### E22.3.1 Selection of relative time series

Let the timing of the break is  $y_0$ . As the break detection of step 22.2 was in annual scale, the month of the first estimated break position is December, by definition.

For auditing possible connected outliers, a relative time series must include the period between  $y_0-1$  and  $y_0+2$ . B6 is applied to select the best relative time series covering the defined period.

### E22.3.2 Flagging outliers

Months of detected outliers are flagged if they have non-interrupted temporal connection with a detected break. Months of outlier periods of  $L < 5$  are considered, while longer outlier periods are excluded. A connection with the break is interrupted when at least 1 non-outlier monthly value with status “observed” separates the break and the outlier. Note that months of status “interpolated” do not produce interruption. The maximal temporal difference between the date of the break and that of a flagged month is 11 months.

### E22.3.3 Audit of flagged values

Flagged months after the break of year  $y_0$  are presented as  $m^\#$  of year  $y_0+1$  and their total number is  $H^\#$ .  $\mathbf{T}$  and  $\mathbf{T}^*$  are used according to their definition in step 21.2. The outliers are retained if (74) is true, while their outlier status is cancelled in the opposite case.

$$\left| \overline{\mathbf{T}_{[y_0-1, y_0]}} - \frac{1}{H^\#} \sum m^\# t_{y_0+1, m^\#} \right| < \left| \overline{\mathbf{T}_{[y_0-1, y_0]}} - \frac{1}{H^\#} \sum m^\# t_{y_0+1, m^\#}^* \right| \quad (74)$$

Similarly, if the flagged months are before the break, the outliers are retained if (75) is true, while their outlier status is cancelled in the opposite case.

$$\left| \overline{\mathbf{T}_{[y_0+1, y_0+2]}} - \frac{1}{H^\#} \sum_{m^\#} t_{y_0, m^\#} \right| < \left| \overline{\mathbf{T}_{[y_0+1, y_0+2]}} - \frac{1}{H^\#} \sum_{m^\#} t_{y_0, m^\#}^* \right| \quad (75)$$

When the status of outlier is cancelled for a month, the observed monthly value (in case of daily homogenization the observed daily values) will be included in  $\mathbf{G}^\#$  and in case of RR homogenization also in  $\mathbf{X}^+$ . The related corrections are made in  $\mathbf{T}$ .

#### E22.4 Monthly precision

For break in year  $y_0$ , the period  $[y_0-1, y_0+2]$  of  $\mathbf{T}$  is examined to find the timing with monthly precision. The break position is expected to be in a narrower, 29-month wide window, i.e. between October of  $y_0-1$  and February of  $y_0+2$ .

##### E22.4.1 Selection of relative time series

The same as step 22.3.1.

##### E22.4.2 Calculation of break position with monthly preciseness

- i) In all homogenization tasks except of sinusoid seasonal cycle  
Optimal step function (see B7) of  $K = 1$  is fitted to the 48-month period of  $\mathbf{T}$ .
- ii) Sinusoid seasonal cycle of inhomogeneities  
Modified step function including sinusoid changes within sections (see B9) of  $K = 1$  is fitted to the 48-month period of  $\mathbf{T}$ .

##### E.22.4.3 Separation of break positions for bi-seasonal detection results

Performed for bi-seasonal RR homogenization.

Let suppose that a common break for the two variables is detected at  $h_0$ . However, a break of rain precipitation does not have sense within the snowy season, and vice versa. Therefore, when  $h_0$  is out of the season for one of the variables, the last month belonging to that season before  $h_0$  will be the break position for that variable.

## E22.5 Merging break lists

Performed for additive variables.

In the homogenization of additive variables, two kinds of break detection are performed. One is the step function fitting in annual scale, while the other is the filtering of outlier periods in monthly scale, as the borders of at least 5-month long outlier periods are considered to be breaks. These together produce two break lists, which are merged here.

Detected breaks are ordered according to dates. Sometimes a break of the same date is detected with both detection methods. Breaks of the two break lists are considered identical when their dates are closer than 5 months. In such cases the date obtained with step function fitting and monthly precision is retained only.

## E22.6 Limitation of the number of synchronous breaks

A similar operation is performed to that of step 10.4. The number of synchronous breaks is kept below the 50% of the time series with homogenized period around the date of the synchronous break with formula (61). However, some details of this step differs from those of step 10.4.

- i) Any two breaks are considered synchronous if their time distance is shorter than 5 months.
- ii) As far as breaks of short-term inhomogeneities occur among the synchronous breaks, breaks detected by step function fitting cannot be removed.
- iii) Between breaks of the same kind detection procedure, their significances are evaluated by Eqs. (62-63) of step 10.4, but in bivariate homogenization the definition of  $Q$  is different (76–78).

Univariate break detection: 
$$Q = (\delta\bar{\mathbf{T}})^2 \quad (76)$$

Sinusoid cycle of inhomogeneities: 
$$Q = (\delta\bar{\mathbf{T}})^2 + (\delta\tilde{\mathbf{T}})^2 \quad (77)$$

Bi-seasonal RR homogenization: 
$$Q = (\delta\overline{\mathbf{T}_{\text{rain}}})^2 + (\delta\overline{\mathbf{T}_{\text{snow}}})^2 \quad (78)$$

- iv) In sinusoid annual cycle of inhomogeneities, the treatment of synchronous breaks is not separated. By contrast, in bi-seasonal RR homogenization the treatment is separated, as the break positions are different for the two variables (step 22.4.3).

## E22.7 Daily precision

This step is performed for additive variables, in daily homogenization.

Data around break position  $y_0$ ,  $m_0$  is examined with the help of 4-month and 12-month wide windows. The default day of a break is the last day of  $m_0$ , and calculations for achieving higher preciseness is provided only when the break has at least comparable size with the standard deviation of daily data in the relative time series (see more details below).

#### E22.7.1 Construction of relative time series of daily data

A 12-month wide symmetric window is edited around the default position of the break. The section between its 5<sup>th</sup> and 8<sup>th</sup> months (both included) is named central section. Type of candidate series:  $\mathbf{G}^+$ , type of reference composites:  $\mathbf{G}^{+*}$ .

Reference composites must cover the defined time window, and they must have at least 0.4 spatial correlation ( $r^*$ ) with the candidate series. Appropriate reference composites are ordered according to  $r^*$ , and maximum 10 of them are included in the construction of the one only reference series ( $\mathbf{F}$ ). The reference composites are weighted by  $r^{*2}$ , then the relative time series ( $\mathbf{T}$ ) is generated by Eq. (10) of B4.

#### E22.7.2 Signal-to-noise ratio

Standard deviation of daily values ( $\sigma^{(d)}$ ) of  $\mathbf{T}$  is calculated for the whole of the 12-month period ( $\sigma_A^{(d)}$ ) and for its central section ( $\sigma_B^{(d)}$ ). The signal-to-noise ratio ( $\alpha^*$ ) is considered sufficient if relation (79) is true.

$$\alpha^* \geq 0.75(\max(\sigma_A^{(d)}, \sigma_B^{(d)}))^2 \quad (79)$$

#### E22.7.3 Calculation of break position with daily preciseness

This step is performed when relation (79) is true.

The central section of  $\mathbf{T}$  is examined. The break position is expected to be maximum 22 days distance from the default day. Optimal step function of  $K = 1$  is fitted (see B7) to the data of the central section.

### 23 Adjustments for inhomogeneities

From this step, the multi-network homogenization is not the same as the 1-network homogenization, as in multi-network homogenization the accuracy of central series is focused.

Time resolution for steps 23.1 – 23.4: monthly in monthly homogenization and in RR homogenization, daily in daily homogenization for additive variables.

### E23.1 Application of ANOVA correction method

Performed in multi-network homogenization. Type of input series:  $\mathbf{G}^\#$ .

The ANOVA correction model (B10) is applied for the variable(s) examined in main step 22. The result will be a vector of adjustment terms ( $\mathbf{Z}^*$ ) (or vectors  $\mathbf{Z}^*$  and  $\mathbf{Z}^{**}$  in bivariate cases) for each time series of the network.

### E23.2 Calculation of adjustment terms backwards from the beginning of the homogenized period

Performed in multi-network homogenization. Section of time series: from the first year of the target period until the last year before the homogenized period.

This step is performed similarly to step 11.2, with two main differences: a) Here the time resolution is monthly or daily; b) the calculation of adjustment terms goes back until the first year of the time series.

Term “long-term adjustment term” ( $z_L$ ) is defined as the minimum of the adjustment term for the first month (first day) of the homogenized period in monthly (daily) homogenization ( $z_A$ ) on one hand, and the average adjustment term for the first 30 years of the homogenized period ( $z_B$ ). Note that in bi-seasonal RR homogenization the first month of the relevant season represents the first month of the homogenized period. With this modified definition (64) is valid for this step. When  $z_L \neq z_A$ ,  $z^*$  gradually changes from  $z_A$  to  $z_L$  during 36 months, going backwards from January of the first year of the homogenized period (in this case it is January also in bi-seasonal RR homogenization) with monthly steps. In daily homogenization, the daily adjustment terms of a given month are uniform before the homogenized period. Note that when the homogenized period is shorter than 30 years,  $z_L = 0$ .

### E23.3 Application of weighted ANOVA

Type of input series:  $\mathbf{G}^\#$ .

#### i) 1-network homogenization

The model described in B11 is applied  $N$  times for a given variable, each time with a different candidate series. The results of the candidate series are retained, and finally  $\mathbf{Z}^*$  ( $\mathbf{Z}^*$  and  $\mathbf{Z}^{**}$  in bivariate cases) is produced for each time series.

#### ii) Multi-network homogenization

The central series is selected to be the candidate series of model B11.  $\mathbf{Z}^*$  ( $\mathbf{Z}^*$  and  $\mathbf{Z}^{**}$  in bivariate cases) of the central series is overwritten with the new results.

#### E23.4 Calculation of adjustment terms backwards from the beginning of the homogenized period

It is the same as step 23.2, but in case of multi-network homogenization, it is applied only to the central series.

#### 23.5 Adjustment terms derived from ensemble results

Section of time series: whole series.

The ensemble cycle has been terminated with the previous step. The mean adjustment terms ( $Z'$ , or  $Z'$  and  $Z''$  in the bivariate cases) are the arithmetic averages of the 9 ensemble results of  $Z^*$  ( $Z^*$  and  $Z^{**}$ ). For the correct adjustments, not only the adjustment terms, but also  $G^{\#}$  ( $X^+$ ) in the homogenization of additive variables (precipitation) are retained from each ensemble cycle, and their arithmetical averages for the 9 ensemble members will be used in the adjustments.

#### 23.6 Adjustment terms for application

Section of time series: whole series including excluded years. Temporal resolution: monthly and daily.

The relations between  $Z$  and  $Z'$  ( $Z$  on the one hand and  $Z'$  and  $Z''$  on the other hand in the bivariate cases) are almost the same as in step 11.3.2 for the treated period, in spite of the temporal resolution of  $Z'$  and  $Z''$  here is monthly or daily.

##### i) Univariate homogenization

$z = z'$  for any month and day, except the cases described in paragraphs v) – viii).

##### ii) Sinusoid cycle of inhomogeneities, monthly resolution

Eq. (66) is valid, with the difference that  $z'$  and  $z''$  hold both annual and monthly indexes.

##### iii) Sinusoid cycle of inhomogeneities, daily resolution

$$z_{y,m,d} = z'_{y,m,d} + 0.55 \sin\left(\frac{2\pi\left(m-3.2+\frac{d}{D_m}\right)}{12}\right) z''_{y,m,d} \quad (80)$$

In (80),  $D_m$  stands for the number of days in month  $m$ .

##### iv) Bi-seasonal RR homogenization

Eqs. (67) and (68) are valid with the difference that  $z'$  and  $z''$  hold both annual and monthly indexes.



v) Missing data in the climate series

Gap-filling was applied only within the treated period, hence data gaps may exist in  $\mathbf{G}^\#$  and  $\mathbf{X}^+$ . For dates without data in series  $\mathbf{G}^\#$  ( $\mathbf{X}^+$  in RR homogenization),  $z = 0$  by definition.

vi) Observed monthly values in excluded years

Let  $y_0$  be an excluded year. The adjustment term is the same as for the first month of the first not excluded year after  $y_0$ .

$$Z_{y_0,m} = Z_{y_0+j,1} \quad (81)$$

In Eq. (81)  $j$  denotes the lowest natural number for which  $y_0+j$  is not an excluded year.

vii) Observed daily values in excluded years

The adjustment term is the same as for the first day of the first not excluded year after  $y_0$  (82).

$$Z_{y_0,m,d} = Z_{y_0+j,1,1} \quad (82)$$

viii) Observed value of excluded year in bi-seasonal RR homogenization (83).

$$Z_{y_0,m,d} = Z_{y_0,m} = Z_{y_0+j,m_1} \quad (83)$$

In Eq. (83)  $m_1$  is the first month of the season including  $m$ .

### 23.7 Execution of adjustments

Section of time series: whole series including excluded years. Time resolution: monthly in monthly homogenization and daily in daily homogenization.

i) Additive variables

$$\mathbf{G}^{*-} = \overline{\mathbf{G}^\#} + \mathbf{Z} \quad (84)$$

In Eq. (84) double stroke denotes ensemble average. With the exception of irregular seasonality of inhomogeneities, this is the final results, i.e.  $\mathbf{G}^{**} \equiv \mathbf{G}^{*-}$ .

In daily homogenization, adjusted values of no outlier filtered series (85) will also be in use.

$$\mathbf{G}^{*-} = \mathbf{G} + \mathbf{Z} \quad (85)$$

With the exception of irregular seasonality,  $\mathbf{G}^* \equiv \mathbf{G}^{*-}$ .

ii) Precipitation (case of daily homogenization) (86).

$$x_{y,m,d}^{**} = \overline{\overline{x_{y,m,d}^+}} e^{z_{y,m,d}} \text{ for every } y,m,d \quad (86)$$

## VIII FINAL OPERATIONS

### 24 Irregular seasonal cycle of inhomogeneities

Concept: The signal-to-noise ratio is less favourable for the assessment of the seasonal cycle of inhomogeneities than for that of the annual mean biases, therefore a part of the seeming seasonal cycles is spurious. Therefore, the method described here reduces the appearance of spurious seasonal changes at the cost of reducing sometimes true amplitudes.

The break detection of steps 22.2 and 22.4 is repeated here with modified parameterization, which is more restrictive towards detecting small or short biases. The break timings are common for any part of the year, but the resulted biases differ. The annual series of monthly values for fixed calendar months are used to assess monthly adjustment terms. The minimum biases of ensemble bias assessments will be used for the calculation of the adjustment terms together with refinements detailed in this main step.

#### 24.1 Selection of the group of relative time series

The relative time series (**T**) constructed by steps 21.1 and 21.2 with the most probable scenario  $p^{(5)}$  are selected.

#### 24.2 Break detection

##### 24.2.1 Selection of relative time series for step function fitting

The same as step 10.1.

##### 24.2.2 Step function fitting

Univariate detection (B7) is applied with  $p_2 = 3.36$ . The minimal distance between two breaks is 5 years.

##### 24.2.3 Selection of relative time series in monthly precision

The same as step 22.3.1.

##### 24.2.4 Monthly precision

The same as the univariate case in step 22.4.

#### 24.2.5 Limitation of the number of synchronous breaks

The same as step 10.4, with the supplement that any two breaks are considered synchronous if their time distance is shorter than 17 months.

### 24.3 Adjustments

#### E24.3.1 ANOVA correction model for annual series of monthly values

Type of input series:  $\mathbf{G}^+$ . Section of time series: homogenized period. Time resolution: annual.

Annual series for the monthly values of fixed calendar months are taken for all the time series in network and for all the 12 months of the year. Model B10 is applied to the set of each calendar month. An ensemble is generated in the way that 1 time series is excluded for each ensemble member, hence a similar ensemble is produced to the one at main step 11 (with the difference that there the ensemble cycle included more operations).

#### E24.3.2 Calculation of adjustment terms backwards from the beginning of the homogenized period

Section of time series: from the first year of the target period until the starting of the homogenized period. Time resolution: annual.

The same as step 11.2, except that the calculations go back until the first year of the time series.

#### 24.3.3 Monthly adjustment terms based on the evaluation of ensemble calculations

Section of time series: whole series. Time resolution: monthly.

The ensemble calculations are evaluated by Eq. (65), in the same way as in step 11.3.1. Once the calculations have been performed for each calendar month in annual resolution, the result adjustment terms (denoted by  $\mathbf{\Omega}^*$ ) are present in monthly resolution for each time series.

#### 24.3.4 Smoothing between adjacent months

Section of time series: whole series. Time resolution: monthly

Here, monthly adjustment terms are indexed by the serial number of month ( $h$ ) from the starting of the time series. A refinement of adjustment terms is provided by the smoothing with (87).

$$\omega_h^\# = 0.3\omega_{h-1}^* + 0.4\omega_h^* + 0.3\omega_{h+1}^* \text{ for all } h \in (2, 3, \dots, H-1) \quad (87)$$

#### 24.3.5 Removal of annual changes resulted by the seasonal adjustments

The previous operations might result in undesired changes in the annual values of the target variable. The seasonal adjustment terms are modified here to remove such annual changes.

i) Within the homogenized period: The mean of the seasonal adjustment terms between adjacent breaks of the step function of step 24.3.1 is subtracted from the adjustment terms (88).

$$\omega'_{y,m} = \omega_{y,m}^\# - \overline{\Omega_k^\#} \text{ for all } y, m \in k, k \in (0, 1, 2, \dots, K) \quad (88)$$

ii) Between the first year of the time series and the starting of the homogenized period: from the seasonal adjustment terms their annual mean is subtracted.

In monthly homogenization  $\Omega \equiv \Omega'$ .

#### 24.3.6 Downscaling to daily adjustment terms

Section of time series: whole series.

Concept: The purpose of Vincent method is to provide linear changes of adjustment terms for days between two adjacent middle-of-months, in a way that monthly adjustment terms remain unchanged. The method exploits the fact that at a simple linear interpolation between middle months (as that is applied at step 11.3.2), the resulted monthly values are determined by the input monthly value of the actual month in 75% and by those of the two adjacent months in 12.5% for each. Introducing auxiliary variable  $\mathbf{A}$ , the monthly values of  $\Omega'$  can be kept during the downscaling (89-91).

$$0.125a_{h-1} + 0.75a_h + 0.125a_{h+1} = \omega'_h \text{ for all } h \in (2, 3, \dots, H-1) \quad (89)$$

$$0.875a_1 + 0.125a_2 = \omega'_1 \quad (90)$$

$$0.125a_{H-1} + 0.875a_H = \omega'_H \quad (91)$$

Note that when the changes between monthly values are nearly linear, or when the accuracy of monthly values is less important, the simple linear interpolation of step 11.3.2. can still be applied. In case of seasonal differences of biases,  $\Omega$  series of daily resolution are provided by these two downscaling methods.

i) In 1-network homogenization

The Vincent method is applied to all series.

ii) In multi-network homogenization

The Vincent method is applied to the central series. Simple linear interpolation (with input monthly values of  $\Omega'$ ) is applied to the other time series.

#### 24.3.7 Adjustment terms in excluded years

Time resolution: monthly and daily.

The example of daily resolution is shown. Let  $y_0$  be an excluded year, and  $j$  is the lowest natural number for which  $y_0+j$  is not an excluded year. Then (92) shows the determination of  $\omega$  for any  $m$  and  $d$  of year  $y_0$ .

$$\omega_{y_0,m,d} = \omega_{y_0+j,m,d} \quad (92)$$

#### 24.3.8 Application of adjustment terms

Section of time series: whole series including excluded years. Time resolution: monthly and daily.

Equations (93-94) are valid in any time resolution.

$$\mathbf{G}^{**} = \mathbf{G}^{*-} + \Omega \quad (93)$$

$$\mathbf{G}^* = \mathbf{G}^{*-} + \Omega \quad (94)$$

Note that at this phase of the procedure,  $\mathbf{G}^*$  is used only in daily homogenization.

### 25 Refinement of outlier periods

Performed in daily homogenization of additive variables.

The subject of the examinations are the outlier periods of 1-4 month length according to the detection results of main step 19. Such periods of any time series are examined one-

by-one here. The goal is to provide daily preciseness for the positions of these outlier periods, and to assess their mean bias from data adjusted in the previous steps. Note that the latter is only for providing information in the output results, since the mean bias of an outlier period is not used for adjustments.

For the examinations, an outlier period is supplied with their adjacent 4-month long periods in its both sides. These supplement sections must not stretch out of the homogenized period, and must not contain detected outliers. When a supplement section does not meet with these conditions, it is shortened as far as the conditions are completed, and in the extreme case its extent can be 0. Time  $d'$  is defined as the distance in days from the starting point of the lengthened period including supplement sections. The pre-estimated dates of the borders of the outlier period ( $d_1^*$  and  $d_2^*$ ) are identical with the first and last days of the detected outlier period in main step 19. Dates  $d_1^*$  and  $d_2^*$  are also the default solution, for cases when the assessment with daily preciseness is denied by the program for the lack of required conditions.

### 25.1 Construction of relative time series

Section of time series: the lengthened outlier period. Type of candidate series:  $\mathbf{Gc}^*$ , type of reference composites:  $\mathbf{G}^{**}$ .

Reference composites must cover the examined period, cannot contain detected outliers within the examined section, and must have at least 0.4 spatial correlation ( $r^*$ ) with the candidate series. The other details are the same as in step 22.7.1.

### 25.2 Positions of outlier periods in daily scale

Section of time series: the lengthened outlier period.

Optimal step function (see B7) of  $K = 2$  is fitted to the data. The final borders of the outlier period ( $d_1'$  and  $d_2'$ ) may have maximum 20 days distance from the default dates ( $d_1^*$  and  $d_2^*$ ). The minimum length of an outlier period is 10 days.

The effect of these results is that a) interpolated values will be provided at main step 26 for all the dates within the outlier period, b) the homogenized values for the dates out of the outlier period will be  $gc + z$  (or  $gc + z + \omega$ , in case of irregular seasonality).

### 25.3 Mean bias of the outlier period

The mean bias ( $q$ ) is the difference between the mean of the values within the outlier period and that within the supplement sections. In (95)  $l_1$  is the length of the outlier period  $[d_1', d_2']$ , and  $l_2$  is the length of the outlier period together with its supplement sections.

$$q = l_1 \overline{\mathbf{G}_{[d'_1+1, d'_2]}^*} - \frac{d'_1 \overline{\mathbf{G}_{[1, d'_1]}^*} + (l_2 - d'_2) \overline{\mathbf{G}_{[d'_2+1, l_2]}^*}}{d'_1 + l_2 - d'_2} \quad (95)$$

#### 25.4 Required conditions

- i) The two supplement sections together must contain at least 3 months of status “observed”, otherwise neither the refinement of the temporal position, nor the assessment of the mean bias will be done. When the refinement on daily scale is denied by the program, the default position of the outlier period will be its final position, and mean bias is not calculated for that.
- ii) At least 3 appropriate reference composites are needed to construct the relative time series at step 25.1. If this condition is not completed, neither the refinement of the temporal position, nor the assessment of the mean bias are done.
- iii) The standard deviation of daily values ( $\sigma^{(d)}$ ) within the outlier period must not be too large in comparison with that in the supplement sections (96).

$$\sigma_{d'_1+1, d'_2}^{(d)} < 2\sigma_{[1, d'_1] \cup [d'_2+1, l_2]}^{(d)} \quad (96)$$

If relation (96) is not completed, the mean bias will not be assessed, but it does not affect the refinement of the temporal position of the outlier period.

#### 26 Infilling data gaps

In 1-network homogenization performed for all series. In multi-network homogenization performed only for the central series. Section of time series: whole series including excluded years. Time resolution: monthly or daily.

Interpolated values based on homogenized observed values of nearby stations will be provided a) for missing data; b) for monthly outlier values of precipitation (in monthly RR homogenization) detected at main step 13; c) for monthly outlier values or outlier periods of maximum 4-month length of additive variables (in monthly homogenization of additive variables), detected at main step 19; d) for outlier periods of daily data (in daily homogenization of additive variables) detected at main step 19 and refined at step 25.2. Spatial correlations ( $r$ ) of step 20.1 will be used here.



## 26.1 Gap filling

B12 – B13 are used, but several other details differ from the previous gap filling routines (e.g. of step 5.2). Let suppose that the candidate series (**Gc**) have a missing value at year  $y_0$ , month  $m_0$ , day  $d_0$ . A potential partner series (**Gs**) must have an observed value at the same date, which must not belong to an outlier period. Beyond this, it must have at least 900 synchronous daily value pairs or 30 synchronous monthly value pairs with **Gc**, in which none of the data can be interpolated or flagged as outlier or as part of a short outlier period. Potential partner series are put into the decreasing order of  $r_{gc,s}$ , and their weights in the interpolation are determined one-by-one using  $r^2$ ,  $p_8$  and  $p_9$ .  $p_8$  is a coefficient modifying the weight ( $w_s$ ) of a partner series relative to  $r^2$ , while  $p_9$  marks the minimum threshold of weights at which partner series are accepted. These parameters are described below in details.

Synchronous value pairs for **Gc** and **Gs** are searched first in year  $y_0$ , then in gradually increasing distances from  $y_0$ . Dates within the same 3-month season are accepted to which  $m_0$  belongs (i.e.  $m \in m_0^*$  for all value pairs). There is no pre-limitation of the window width around  $y_0$ , but the collection of value pairs terminates when the number of daily value pairs reaches 1800 or the number of monthly value pairs reaches 60, hence in nearly complete datasets the window width is 20-25 years. Let  $\lambda$  stand for the mean distance of value pairs from  $y_0$  in years, then the weights of the partner series are calculated by (97-98).

$$p_8 = 1 - 0.1\ln\lambda + 0.1672 \quad (97)$$

$$w'_s = p_8 r_{gc,s}^2 \quad (98)$$

If the homogenized period of **Gs** includes the date  $d_0-m_0-y_0$ , then the final weight is  $w_s = w'_s$ , if the homogenized period does not include it, but the treated period yes, then  $w_s = 0.5w'_s$ , while if it is out of the treated period of **Gs**, then  $w_s = 0.33w'_s$ . Note that when both **Gc** and **Gs** are complete in the 21-year symmetric window around month  $m_0-y_0$ , then  $\lambda = 5.32$  and  $p_8 = 1$ . When  $\lambda = 10$  or  $\lambda = 23.9$  or  $\lambda = 64.9$ , then  $p_8$  is 0.937, 0.85, 0.75, respectively.

When the weights have been calculated for at least 4 partner series, the weights are ordered, and the minimum threshold of weights is set by (99).

$$p_9 = \frac{w^{(1)} + w^{(2)} + w^{(3)}}{3} - 0.2 \quad (99)$$

In (99), the upper index of  $w$  denotes the serial number of weight in the rank order. Partner series with  $w < p_9$  are excluded, and in any case the maximum number of partner series is 10.

The calculations for Eqs. (97-99) are repeated for each potential partner series participating in the interpolation for date  $d_0-m_0-y_0$  of **Gc**. The procedure of selecting partner series and determining their weights in the interpolation is repeated for each missing data of the dataset.

## 26.2 Giving back mean seasonal cycle

Performed for additive variables.

Climatic mean seasonal values defined by Eq. (9) are added to the homogenized deseasonalised series (100).

$$x_{y,m,d}^{**} = g_{y,m,d}^{**} + \overline{X_{s,m}} \quad (100)$$

Eq. (100) is performed for all dates  $(y,m,d)$  where the homogenized series have either an observed value or an interpolated value. Note that in RR homogenization the results of (86) include the seasonality.

## 26.3 Reliability indicators

Reliability indicators do not have role in the accuracy of the homogenization, but they provide information about the data. The Manual describes the reliability indicators, except that no detailed information is given there about codes 3...7.

Each of codes 3...7 indicates interpolated data, but lower codes indicate higher quality (i.e. interpolation with more or better correlating station series). The number of partner series ( $N''$ ) and the total weight of partner series ( $W$ ) determine the codes.

Code 1 – Homogenized observed data

Code 2 – Observed data within the treated period, but out of the homogenized period

Code 3 – Interpolated data,  $W \geq 3$

Code 4 – Interpolated data,  $2 \leq W < 3$

Code 5 – Interpolated data,  $N'' > 2$  and  $1 \leq W < 2$

Code 6 – Interpolated data, ( $N'' = 2$  and  $W \geq 0.3$ ) or ( $N'' > 1$  and  $0.3 \leq W < 1$ )

Code 7 – Interpolated data,  $N'' = 1$  or ( $N'' > 1$  and  $W < 0.3$ )

Code 8 – long-term climatic mean value, as spatial interpolation is not possible

Code 9 – Missing data without gap filling (it may occur when the user does not require gap filling)

Code 0 – Observed data either out of the treated period or in an excluded year

## 27 Elimination of physical outliers

In 1-network homogenization performed for all series. In multi-network homogenization performed only for the central series. Section of time series: whole series including excluded years. Time resolution: monthly and daily.

When values near to the threshold(s) of the range of a climatic element frequently occur, the occurrence of physical outliers in the adjusted values is not very rare. For instance, a daily SS can be 0, or a daily mean HH can be 100 (%), and after the

adjustments the new values might fall out of the physically possible ranges. Note that (i) the occurrence of physical outlier monthly values is much rarer than that of the daily values; (ii) the wrong manual definition of outlier thresholds might result in the exclusion of true climatic values as physical outliers.

### 27.1 Adjustments for physical outliers

When a value falls out of its defined physical range  $[x_{\min}, x_{\max}]$ , that value is substituted with the closest physically acceptable value. For instance, if at day  $d^*$   $HH_{d^*} = 102$ , it is transformed to  $HH_{d^*} = x_{\max}(HH) = 100$ .

### 27.2 Adjustments for keeping monthly values unchanged

Performed in daily homogenization.

As in other parts of the homogenization procedure, daily adjustments are not allowed to change the monthly values. Let  $b$  be the monthly bias of  $x_{y,m}$  caused by the adjustments of daily outliers in step 27.1. Before the adjustments will be done at this step,  $\Delta x_{y,m} = b$ . A step-by-step adjustment of daily values starts here to reduce  $\Delta x_{y,m}$ . For  $b > 0$ , the individual adjustments in the homogenization of RR, SS or RS are shown by (101), while those for the homogenization of other variables are shown by (102).

$$x_{y,m,d}^{(i+1)} = \max(x_{\min}, x_{y,m,d}^{(i)} - 0.01b) \quad (101)$$

$$x_{y,m,d}^{(i+1)} = \max(x_{\min}, x_{y,m,d}^{(i)} - 0.01bD_m) \quad (102)$$

The upper index of  $x$  denotes the serial number of the iteration. The adjustments of Eqs. (101-102) are performed from the first day until the last day of month  $m$  of year  $y$  as many times as they are necessary, and they are finished at any day if  $|\Delta x_{y,m}| < 0.01b$ . When  $b < 0$ , the daily values are raised by the same gradualness as the reductions are done for  $b > 0$ .

### 27.3 Calculation of monthly values

Performed in daily homogenization.

The monthly values are the arithmetical averages of the daily values in the homogenization of TT, HH, VV, PP or SP. The monthly values are the sums of the daily values in the homogenization of RR, SS or RS.

In a few output items, interpolated values are not included for substituting missing data. In such cases, the status of monthly data is missing, and missing data code

is applied, if the number of missing daily data within a given month is higher than 0 (higher than 7) in RR homogenization (in the homogenization of other variables than RR). If a homogenized monthly value is calculated from less observed daily values than the number of days in month, its correct consideration is straightforward when the monthly value is arithmetical average. When a monthly total must be calculated from an incomplete set of daily values, the unbiased solution is shown by (103).

$$x_{y,m}^{**} = \frac{D_m}{D'_{y,m}} \sum_Y x_{y,m,d}^{**} \quad (103)$$

In Eq. (103),  $Y$  stands for the cluster of days with observed values in month  $m$  of year  $y$ , while  $D'$  presents the number of days with observed data in cluster  $Y$ .

## IX LITERATURE

### L1 Description of earlier ACMANT versions and sources of ACMANT

Caussinus, H. and Lyazrhi, F. (1997) Choosing a linear model with a random number of change-points and outliers. *Ann. Inst. Statist. Math.* 49(4):761-775.

Caussinus, H. and Mestre, O. (2004) Detection and correction of artificial shifts in climate series. *J. Roy. Stat. Soc. C* 53:405-425. <http://doi.org/10.1111/j.1467-9876.2004.05155.x>

Domonkos, P. (2011) Adapted Caussinus-Mestre Algorithm for Networks of Temperature series (ACMANT). *Int. J. Geosci.* 2:293-309. <http://doi.org/10.4236/ijg.2011.23032>.

Domonkos, P. (2014) The ACMANT2 software package. In: Eighth Seminar for Homogenization and Quality Control in Climatological Databases and Third Conference on Spatial Interpolation Techniques In Climatology and Meteorology (ed. Lakatos, M., Szentimrey, T. and Marton A.), WMO WCDMP-84, Geneva, Switzerland, 46-72.

Domonkos, P. (2015) Homogenization of precipitation time series with ACMANT. *Theor. Appl. Climatol.* 122:303-314. <http://doi.org/10.1007/s00704-014-1298-5>.

Domonkos, P. and Coll, J. (2017) Homogenisation of temperature and precipitation time series with ACMANT3: Method description and efficiency tests. *Int. J. Climatol.* 37:1910-1921. <http://doi.org/10.1002/joc.4822>

Hawkins, D.M. (1972) On the choice of segments in piecewise approximation. *J. Inst. Math. Appl.* 9:250–256. <http://doi.org/10.1093/imamat/9.2.250>.

Lindau, R. and Venema, V. (2018) On the reduction of trend errors by the ANOVA joint correction scheme used in homogenization of climate station records. *Int. J. Climatol.* 38:5255-5271. <http://doi.org/10.1002/joc.5728>

Peterson, T.C. and Easterling, D.R. (1994) Creation of homogeneous composite climatological reference series. *Int. J. Climatol.* 14:671–679.

Szentimrey, T. (2010) Methodological questions of series comparison. In: 6th Seminar for Homogenization and Quality Control in Climatological Databases (Ed. Lakatos, M., Szentimrey, T., Bihari, Z. and Szalai, S.) WMO WCDMP-76:1-7.

Vincent, L.A., Zhang, X., Bonsal, B.R. and Hogg, W.D. (2002) Homogenization of daily temperatures over Canada. *J. Clim.* 15:1322–1334. [http://doi.org/10.1175/1520-0442\(2002\)015<1322:HODTOC>2.0.CO;2](http://doi.org/10.1175/1520-0442(2002)015<1322:HODTOC>2.0.CO;2)

## **L2 Properties of ACMANT or its routines**

Domonkos, P. and Coll, J. (2017) Time series homogenisation of large observational datasets: The impact of the number of partner series on the efficiency. *Clim. Res.* 74:31-42. <http://doi.org/10.3354/cr01488>.

Domonkos, P. and Coll, J. (2019) Impact of missing data on the efficiency of homogenization: Experiments with ACMANTv3. *Theor. Appl. Climatol.* 136:287-299. <http://doi.org/10.1007/s00704-018-2488-3>.

Lindau, R. and Venema, V. (2013) On the multiple breakpoint problem and the number of significant breaks in homogenization of climate records. *Időjárás Q. J. Hungarian Meteorological Service*, 117:1–34.

Lindau, R. and Venema, V. (2018) On the reduction of trend errors by the ANOVA joint correction scheme used in homogenization of climate station records. *Int. J. Climatol.* 38:5255–5271. <http://doi.org/10.1002/joc.5728>.

Lindau, R. and Venema, V. (2018) The joint influence of break and noise variance on the break detection capability in time series homogenization. *Adv. Stat. Clim. Meteorol. Oceanogr.* 4:1–18. <http://doi.org/10.5194/ascmo-4-1-2018>

Lindau, R. and Venema, V. (2019) A new method to study inhomogeneities in climate records: Brownian motion or random deviations? *Int. J. Climatol.* 39:4769–4783. <http://doi.org/10.1002/joc.6105>

## **L3 ACMANT or its routines in method comparison studies**

Domonkos, P. (2011) Efficiency evaluation for detecting inhomogeneities by objective homogenisation methods. *Theor. Appl. Climatol.* 105:455–467. <http://doi.org/10.1007/s00704-011-0399-7>

Domonkos, P. (2013) Efficiencies of inhomogeneity-detection algorithms: comparison of different detection methods and efficiency measures. *J. Climatol.* pp15. <http://doi.org/10.1155/2013/390945>

Domonkos, P., Venema, V. and Mestre, O. (2011) Efficiencies of homogenisation methods: our present knowledge and its limitation. *Seventh Seminar for Homogenisation and Quality Control in Climatological Databases* (ed. Lakatos, M., Szentimrey, T. and Vincze, E.), WMO-WCDMP-78, Geneva, Switzerland, 19–32.

Guijarro, J.A., López, J.A., Aguilar, E., Domonkos, P., Venema, V., Sigró, J. and Brunet, M. (2017) Comparison of homogenization packages applied to monthly series of temperature and precipitation: the MULTITEST project. *Ninth Seminar for Homogenization and Quality Control in Climatological Databases and Fourth Conference on Spatial Interpolation Techniques In Climatology and Meteorology* (ed. Szentimrey, T., Lakatos, M. and Hoffmann, L.), WMO WCDMP-85, Geneva, Switzerland, 46–62.

Killick, R.E. (2016) Benchmarking the Performance of Homogenisation Algorithms on Daily Temperature Data. PhD thesis, University of Exeter, UK.  
<https://ore.exeter.ac.uk/repository/handle/10871/23095>

Menne, M.J. and Williams, C.N. Jr. (2005) Detection of undocumented changepoints using multiple test statistics and composite reference series. *J. Clim.* 18:4271–4286.  
<http://doi.org/10.1175/JCLI3524.1>

Venema, V., Mestre, O., Aguilar, E., Auer, I., Guijarro, J.A., Domonkos, P., Vertacnik, G., Szentimrey, T., Štěpánek, P., Zahradnicek, P., Viarre, J., Müller-Westermeier, G., Lakatos, M., Williams, C.N., Menne, M.J., Lindau, R., Rasol, D., Rustemeier, E., Kolokythas, K., Marinova, T., Andresen, L., Acquaotta, F., Fratianni, S., Cheval, S., Klancar, M., Brunetti, M., Gruber, C., Duran, M.P., Likso, T., Esteban, P. and Brandsma, T. (2012) Benchmarking monthly homogenization algorithms. *Clim. Past*, 8:89–115. <http://doi.org/10.5194/cp-8-89-2012>