RESEARCH ARTICLE



Uncertainty evaluation of Climatol's adjustment algorithm applied to daily air temperature time series

Oleg Skrynyk^{1,2} | Enric Aguilar¹ | Jose Guijarro³ | Luc Yannick Andreas Randriamarolaza^{1,4} | Sergiy Bubin⁵

¹Center for Climate Change, C3, Geography Department, Universitat Rovira i Virgili, Tarragona, Spain

²Division of Atmospheric Physics, Ukrainian Hydrometeorological Institute, Kyiv, Ukraine

³State Meteorological Agency (AEMET), Balearic Islands Office, Palma de Mallorca, Spain

⁴Division de la Météorologie appliquée, Direction Générale de la Météorologie, Madagascar

⁵Department of Physics, Nazarbayev University, Nur-Sultan, Kazakhstan

Correspondence

Oleg Skrynyk, Center for Climate Change, C3, Geography Department, Universitat Rovira i Virgili, 15 C. Joanot Martorell, Vila-seca, 43480, Tarragona, Spain. Email: oleg.skrynyk@urv.cat

Funding information

The work was performed in the frame of the INDECIS project, funded by FORMAS (SE), DLR (DE), BMWFW (AT), IFD (DK), MINECO (ES), ANR (FR) with co-funding by the European Union, Grant/Award Number: 690462; Ministry of Education and Science of Kazakhstan, Grant/Award Number: BR05236454; Nazarbayev University, Grant/Award Number: 090118FD5345

Abstract

The present study investigated the uncertainty associated with Climatol's adjustment algorithm applied to daily minimum and maximum air temperature. The uncertainty quantification was performed based on several numerical experiments and the benchmark data that were created in the framework of the INDECIS project. Using a complex approach, the uncertainty was evaluated at different levels of detail (day-to-day evaluation through formalism of random functions and six statistical metrics) and time resolution (daily and yearly). However, only the main source of potential residual errors was considered, namely station signals introduced into a raw data set to be homogenized/ adjusted. Other influencing factors, such as the averaged correlation between a candidate and references, were removed from the analysis or kept almost unchanged. According to our calculations, the Climatol's adjustment uncertainty, evaluated on the daily scale, varies over time. The width of the residual errors distribution in summer months is substantially less compared with wintertime. The slight seasonality is also observed in the means of the residual errors. The uncertainty evaluation based on the statistical metrics usually neglect such seasonal non-stationarity of the residual errors providing only assessments averaged over time. On the other hand, the metrics provide detailed information regarding both types of the residual errors, systematic and scatter. The metrics values confirmed good capability of the Climatol software to remove the systematic errors related to jumps in the means, while the scatter errors are removed from the raw time series less efficiently. On the yearly scale, the uncertainty evaluation was performed for the yearly temperature data and several climate extreme indices. Both types of errors are removed well in the yearly time series of the air temperature and the extreme indices. The metrics values showed a significant reduction of the Climatol's adjustment uncertainty. A substantial decrease of the linear trend errors in the yearly time series can also be observed.

K E Y W O R D S

Climatol, homogenization adjustment, INDECIS, minimum and maximum daily air temperature, uncertainty

nal 🛛 📉 RMetS

1 | INTRODUCTION

The detection of modern climate change and analysis of climate variability and extreme events on the national, regional, or even global scales are mainly performed based on a statistical analysis of time series of measured meteorological variables such as air temperature and precipitation (e.g., Alexander et al., 2006; Klein Tank et al., 2009; Hartmann et al., 2013). However, in order to extract accurate and reliable conclusions from the analysis it is necessary first to homogenize raw data sets usually artefacts containing many spurious (inhomogeneities) (Aguilar et al., 2003; Trewin, 2010). By performing a homogenization, one aims to remove the inhomogeneities (abrupt shifts/jumps, gradual trends, outliers, etc.) and approximate the data to the real climate signal, that took place in some area. Usually the homogenization procedure allows to improve the consistency of the data, which can be seen in the process of a statistical comparison of the raw and homogenized time series (e.g., Mamara et al., 2014; Prohom et al., 2016; Osadchyi et al., 2018; Yosef et al., 2018; Fioravanti et al., 2019; Skrynyk et al., 2019; Dumitrescu et al., 2020). However, the question that may remain unclear is: how far are the homogenized data from the true climate signal? Or, in other words, what potential uncertainties could still be present in the data homogenized by means of some homogenization algorithm or software? It is a very important yet largely overlooked issue, because the climate signal (clean data) is essentially unknown and it is impossible to conduct a direct quantitative comparison and evaluation of the homogenization results. At the same time, understanding the uncertainties and their causes is vital for the correct interpretation of outputs of any predicting model (e.g., Iman and Helton, 1988), including homogenization software.

The problem of climate data homogenization can be divided into two sub-problems, namely the detection of discontinuities (most probable dates of potential inhomogeneities) and adjustment of inhomogeneous data (some segments of raw time series) to homogeneous state. Both sub-problems might produce a certain number of common errors, which deviate the homogenized data from the true climate signal. Evaluation of the efficiency of the detection algorithms has been performed in many works (e.g., Ducré-Robitaille et al., 2003; DeGaetano, 2006; Reeves et al., 2007; Domonkos, 2011; Kuglitsch et al., 2012; Venema et al., 2012; Willett et al., 2014; Killick, 2016; Yozgatligil and Yazici, 2016; Coll et al., 2020). Assessment of the performance of the adjustment methods has also been considered (e.g., Della-Marta and Wanner, 2006; Mestre et al., 2011; Trewin, 2013; Squintu et al., 2020). In both cases, the evaluation was mainly performed in a relative form, that is, several homogenization algorithms are usually compared in order to define which one gives the best output and is most suitable for practical applications. Such relative comparison is usually performed based on some benchmark data. However, the quantification of the uncertainties of the homogenization procedures has been considered in few works (e.g., Lindau and Venema, 2016; Trewin, 2018; Vincent et al., 2018). Lindau and Venema (2016) studied the uncertainty of the multiple breakpoint detection algorithms applied to the yearly climate time series. To do so, they defined a probability distribution for possible shifts of the detected break from its true position based on a theoretical approach. According to their findings, the probability of the shifts or, in other words, detection errors, can be described statistically by a Brownian motion with drift. Vincent et al. (2018) and Trewin (2018) evaluated the uncertainty of the homogenization adjustment algorithms applied to the daily air temperature time series. In both works, parallel measurements of temperature were used in order to assess potential residual errors. However, the uncertainty of the adjustment was quantified using different methodology. In (Vincent et al., 2018) the remaining errors in corrected time series were evaluated through two statistical metrics, the root mean square error (RMSE) and the percentage of days within 0.5°C (POD05) that were calculated based on daily data. As described in the paper, RMSE and POD05 were used to assess the uncertainty in the mean and extreme temperature values, respectively. In (Trewin, 2018) the uncertainty is also evaluated through some statistical indicators, but they were calculated on the seasonal and annual scales. The uncertainty was defined as the standard deviation (SD) of the indicator values that were obtained by repeating calculations for slightly different adjustment conditions (changing the set of reference stations, their number, etc.). It is important to note that in spite of intuitively clear meaning of the term 'uncertainty', which can be simply interpreted as a range or a distribution of possible residual errors, there is no unique methodology how it can be quantified for the homogenization/adjustment of climate data.

The objective of this work is to evaluate the uncertainty associated to the adjustment of the daily maximum and minimum temperature series using Climatol (Guijarro, 2018; http://climatol.eu/). In order to focus on Climatol's adjustment algorithm, we base our analysis on the assumption of perfect detection. It should be emphasized that the problem of the uncertainty evaluation of the homogenization adjustment is particularly important when dealing with daily time series, since the climate data of such time resolution is the basis for many modern climatological studies (e.g., detection, monitoring, and attribution of changes in climate extremes). In order to achieve our goal we used benchmark data sets (Aguilar *et al.*, 2018; Pérez-Zanón *et al.*, 2018) generated in the framework of the European project INDECIS (Integrated approach for the development across Europe of user oriented climate indicators for GFCS high-priority sectors: agriculture, disaster risk reduction, energy, health, water, and tourism) (INDECIS, 2018).

The methodology proposed in the present work and applied to Climatol can be generalized for other homogenization software, which is capable of adjusting daily time series of climatological variables in automatic mode with predefined break points. We believe our findings should be useful for developers of homogenization methods and software as well as for the end users of that software as they provide an insight of what kind of errors they might expect after applying the homogenization adjustment.

2 | DATA AND METHODS

2.1 | The Climatol homogenization software

The R package Climatol is a homogenization software that has been widely used in recent years for removing inhomogeneities from collections of raw time series of different climate variables and different time resolution (e.g., Mamara et al., 2013; Sanchez-Lorenzo et al., 2015; Guijarro et al., 2018; Meseguer-Ruiz et al., 2018; Azorin-Molina et al., 2019; Coll et al., 2020; Dumitrescu et al., 2020). The effectiveness of the software has been evaluated in several benchmark tests (Venema et al., 2012; Killick, 2016; Guijarro et al., 2017; Guijarro et al., 2019) where it demonstrated good results, which are comparable in terms of accuracy to other well established and tested homogenization algorithms. According to the benchmarks, both part of the homogenization procedure in Climatol, namely the detection and adjustment, work well and allow to remove different types of artefacts (abrupt shifts/jumps, gradual trends, outliers, etc.), thereby increasing the consistency of raw data sets. One of Climatol's feature characteristics is that it can be used automatically, which significantly increases its applicability to large data sets such as the European Climate Assessment and Dataset (ECA&D) (Klein Tank et al., 2002). Several versions of the software have been released since its creation. In our work, we used Climatol 3.1.1., available through CRAN (https:// cran.r-project.org/package=climatol).

The Climatol detection method (Guijarro, 2018) is based on the standard normal homogeneity test (SNHT) (Alexandersson, 1986; Alexandersson and Moberg, 1997). For any candidate time series, Climatol uses data from neighbouring stations to create a single composite reference series as their optionally weighted average.

Climatol first normalizes the data and infills the original missing values (as well as those that were generated by the outlier deletion or the splitting operations on the later computing steps) through an iterative process in which the main statistical properties of time series, namely means and SD, are recalculated at every iteration until their stationary values are obtained. Once the means become stable, all data are normalized and estimated (whether existing or missing, in all of the series) by means of respective value from the composite reference series, that is, as a weighted average of a prescribed number of the nearest available data. From the statistical point of view, the approach used is equivalent to applying a type II linear regression model (Sokal and Rohlf, 1969), which is justified because all climatic time series from a network of stations under study usually have similar errors. At the next step, the normalized original data and their estimates are used to create time series of anomalies (the estimated values are subtracted from the observed ones), which in turn are exploited to find and eliminate outliers and to detect inhomogeneities by applying SNHT. Since SNHT is a test originally devised for finding a single break point in a series, it is applied iteratively, splitting the candidate time series or its segment into two parts every cycle until no inhomogeneous segments are found. Moreover, during iterations, the test is applied twice: (a) to stepped overlapping temporal windows and after that (b) to complete series. Such a two-stage procedure allows to minimize detection errors that occur when two or more shifts in the mean of similar size could mask its results. Finally, all homogeneous sub-periods/segments originate a corresponding number of new series spanning the entire period of study which are reconstructed by using new estimated values to fill in all missing data/segments.

2.2 | The INDECIS benchmark data sets

In the scope of the INDECIS project (see www.indecis. eu), two different collections of benchmark time series were created, which cover two regions in Europe with different climate, namely southern Sweden and Slovenia (Aguilar *et al.*, 2018; Pérez-Zanón *et al.*, 2018). Each collection contains the daily series of nine essential climate variables (cloud cover, wind speed, relative humidity, sea level pressure, precipitation amount, snow depth, sunshine duration, maximum and minimum air temperature) over the period of 1950–2005. Each benchmark data set consists of clean data, extracted from the output of the Royal Netherlands Meteorological Institute (KNMI) Regional Atmospheric Climate Model (RACMO) version 2, driven by Hadley Global Environment Model 2—Earth System (MOHC-HadGEM2-ES) (Collins *et al.*, 2008), and inhomogeneous data, created by introducing realistic breaks and errors. Missing values and other quality problems (different from biases) were also added to generate other flavours of the perturbed benchmarks, however, they were not used in our study. The RACMO model was chosen due to its high spatial resolution $(0.11^{\circ} \times 0.11^{\circ})$ and the daily time step of the output provided: gridded time series of essential climate variables.

In our study, we used only the maximum (TX) and minimum (TN) air temperature benchmark data sets for the southern Sweden (Figure 1a). Both data sets contain 100 'stations', a subset of the RACMO grid points chosen to imitate stations spatial distribution. Their geographical locations on the domain under study are shown in Figure 1b.

The introduction of biases in the homogeneous series was done by simulating relocations. First, the closest pairs of the RACMO grid time series were used to build a database of differences between nearby locations. Then, for every random sub-period to perturb in the homogeneous series, a difference was randomly chosen, modified by a random factor drawn from N(2, 0.2) to enhance the lower variability of modelled series, and applied to bias the sub-period. The total numbers of break points introduced into TN and TX clean time series are 258 and 280, respectively. That is, the mean break frequency was set to ~4/~5 (TN/TX) in 100 years, as it was found in previous studies on European series Domonkos, Venema 2011; (e.g., et al., 2012; Domonkos, 2017). Figure 2 represents the time distribution of the break points, while Figure 3 shows the distribution of the number of stations/time series with respect to the number of breaks in one time series.

Due to the daily time resolution and the way that was used to create the realistic, to the extent possible, station signals (considered here as the time series of the introduced errors, see an example in Figure 7a below), they are characterized by intensive noise presence at each of the homogeneous segments except for the last ones. That makes it difficult to precisely define the factors and amplitudes of the shifts at the break points. Nevertheless, we estimated such parameters by averaging the corresponding sub-periods of the error time series. Thus, in our case the factors are mean values of errors at the homogeneous segments, while the amplitudes are differences between pairs of two consecutive factors: between the means at previous and next segments. As can be seen from Figure 4, where the histograms of the factors and amplitudes are presented, their range for TN, approximately from -6 to 6° C (Figure 4a,c), is wider comparing to TX, (-3; 3) (°C) (Figure 4b,d). This was deliberately introduced into the benchmark to mimic real effects such as those related to larger local microclimate differences at nights comparing to daylight period of days (e.g., Brunet et al., 2008). Beside the factors and amplitudes, the homogeneous segments can also be characterized by SD of errors. Figure 5 shows their histograms for TN and TX time series. The mean and SD of the errors on the homogeneous segments can be combined in a single parameter



FIGURE1 (a) The domain of the southern Sweden (inside of the red rectangular frame) and (b) locations of the 'stations' (the subset of the RACMO grid points, shown as black dots) on it



FIGURE 2 Number of break points per year introduced to clean (a) TN and (b) TX air temperature time series



FIGURE 3 Distribution of the number of stations/time series with respect to the number of break points in one time series: (a) TN, (b) TX

called as signal to noise ratio. But in our work, we consider them separately. It is worth noting that similar to real relocation the introduced errors influence mainly on the average daily temperature causing jumps in the means. However, the extreme daily temperature should be also affected.

The presented statistical properties of the break points and respective homogeneous segments in the introduced station signals are close to reality. Such conclusion is supported by many homogenization results of real data sets where similar statistical features of inhomogeneities have been found (e.g., Brunet *et al.*, 2008; Trewin, 2018).

2.3 | Methodology used to evaluate uncertainty of homogenization adjustment

In order to describe our approach to the evaluation of Climatol's adjustment uncertainty, we first introduce the formalism and present some graphical illustrations. Let

$$\boldsymbol{X}^{I}, \, \boldsymbol{X}^{H}, \, \text{and} \, \boldsymbol{X}^{C}$$
 (1)

be inhomogeneous, homogenized, and clean daily data, respectively. X^{I} and X^{C} can be also referred to as

raw and homogeneous data, correspondingly. All these data sets are collections of time series

$$\boldsymbol{X} = \{x_{ij}\}, \ i = 1, ..., \ M, \ j = 1, ..., N,$$
(2)

where M is the number of meteorological stations considered and N is the number of time steps/days. From the mathematical point of view, X is a rectangular matrix with dimension of $M \times N$. Let X_k , which is the *k*-th row in (2), denote the entire time series for the *k*-th station. The homogenization adjustment can be formally thought as mapping g that transforms the input matrix X^I in to the output one X^{H}

$$\boldsymbol{X}^{I} \stackrel{g}{\to} \boldsymbol{X}^{H}.$$
 (3)

 X^{C} is the reference, etalon result for the outputs.

Based on the data available in (1), time series of real, E^{R} , detected, E^{D} , and homogenization, E^{H} , errors can be calculated:

$$\boldsymbol{E}^{R} = \boldsymbol{X}^{I} - \boldsymbol{X}^{C}, \ \boldsymbol{E}^{D} = \boldsymbol{X}^{I} - \boldsymbol{X}^{H}, \ \boldsymbol{E}^{H} = \boldsymbol{X}^{H} - \boldsymbol{X}^{C}.$$
(4)

Specifically in our case, E^R is a collection of station signals (or, more precisely, station signals plus noise; but we will call them as station signals for simplicity) that were introduced into the clean data X^C . E^H is a data set of residual errors that might be still present in the homogenized or adjusted series X^H . The error datasets E^R , E^D and E^H are also $M \times N$ -matrices: $E = \{e_{ij}\}, i = 1, ..., M, j = 1, ..., N$.

Figure 6 shows some typical examples of the time series associated with the same (k-th) station. They were extracted from the TN raw, homogenized by means of the Climatol software, and clean data sets. Figure 7 shows the corresponding error time series (4), calculated from the data given in Figure 6. All figures can be also interpreted as graphical representations of the k-th rows in the respective matrices. We will refer to both figures throughout this paper to illustrate the configuration and layout of our numerical experiments and results.

The main object of our study is the matrix E^{H} : we want to know how large could be the residual errors in the adjusted data or, in other words, how large could be the departure of the adjustment prediction X^{H} from the reference, etalon result X^{C} . According to, for example, Walker *et al.* (2003), such departure is usually called 'uncertainty'. Typically, there exist multiple reasons, referred to as sources of the uncertainty (Jakeman *et al.*, 2006), which may affect the adjustment performance and magnitude of the errors in E^{H} . Therefore, in

order to evaluate the uncertainty of the homogenization adjustment we must consider all these sources—the whole credible range of every uncertain input and parameter of the adjustment software—and define the effective width of the corresponding probability distribution of the residual errors (Domonkos and Efthymiadis, 2013). The wider the error distribution, the more uncertain the software prediction X^{H} is.

The residual errors of the homogenization adjustment E^{H} should depend on the introduced errors E^{R} . The more complex station signals in E^{R} (e.g., the larger number of break points, the higher amplitudes of shifts, etc.), the larger residual errors should be expected. Thus, to clarify how wide the distribution of the potential remaining errors could be, we need to consider a large number of different yet real variants of E^{R} . Performing the homogenization adjustment for each of them provides a respective ensemble of Climatol's outputs, necessary for the uncertainty quantification.

The result of the homogenization adjustment should also depend on other factors, such as the mean correlation between candidate and reference time series (Szentimrey, 2008; Guijarro, 2011; Domonkos and Coll, 2017), number of reference series (Trewin, 2018), and so forth. However, in the present study we focus only on the influence of the station signals on the adjustment result. That is, we aim to quantify the adjustment uncertainty, which comes from a single source only: the errors introduced into the input data to be adjusted. The sensitivity of Climatol's adjustment to other possible factors will be addressed in our future works.

2.3.1 | The concept of a random field/ function applied to the residual errors E^H

The considerations presented above suggest an appropriate theoretical model for E^H that can provide a basis for further calculations and can make calculation results more solid, both statistically and theoretically. Since we are going to consider an ensemble of different realizations of E^H , it is natural to assume that E^H is a random field or, more generally, a random function, that is given at the limited number $(M \times N)$ of discrete points in space and time domains, D and T, respectively. Therefore, in order to evaluate the homogenization adjustment and to quantify the adjustment uncertainty we have to define and study statistical properties of the random field E^H . According to the theory, a multidimensional $(M \times N)$ dimensional) probability distribution function

$$f_{M \times N} \left(e_{11}^{H}, e_{12}^{H}, ..., e_{1N}^{H}, e_{21}^{H}, ..., e_{2N}^{H}, ..., e_{MN}^{H} \right)$$
(5)



FIGURE 4 Histograms of the factors (a, b) and amplitudes (c, d) of the shifts at break points, that were introduced to TN (a, c) and TX (b, d) clean data sets. The frequency/count was normalized by the total number of the breaks. The factors/amplitudes were estimated by averaging homogeneous segments in the time series of the introduced error



FIGURE 5 Histograms of *SD* of the introduced errors at the homogeneous segments: (a) TN, (b) TX. The frequency/count was normalized by the total number of the breaks



FIGURE 6 Examples of TN time series belonging to the same (*k*-th) station extracted from the inhomogeneous X^{I} (a), homogenized X^{H} (b) and clean X^{C} (c) data sets

provides complete and the most detailed description of E^{H} . Based on $f_{M \times N}$ it is possible to derive multidimensional probability distribution of the residual errors in any of M meteorological stations. For instance, for k-th station we get $f_N(e_{k1}^H, e_{k2}^H, ..., e_{kN}^H)$. The f_N is obtained by integrating $f_{M \times N}$ with respect to its all arguments except $e_{k1}^H, e_{k2}^H, ..., e_{kN}^H$. Function $f_1(e_{kl}^H)$ defines probability distribution of the residual error in k-th meteorological station (i = k) and l-th day (j = l).

In the most general case, a random field might be non-stationary in time and heterogeneous in space. In this situation, the simplest statistical properties of the random field defined in a single point of the space-time domain, such as the mean or *SD*, vary in the domain. On the contrary, when the field is stationary and homogeneous, these statistical moments are constant in time and space. Specifically to the homogenization adjustment, we can expect E^H to be non-stationary (e.g., due to seasonal cycle in temperature time series) and heterogeneous (e.g., due to possible different topography in *D* and, as a result, different local correlation between temperature time series). Such peculiarities of E^{H} , namely non-stationarity and spatial heterogeneity, make its analysis more difficult. In particular, that means we cannot use the ergodic assumption in order to calculate statistical properties of E^{H} based on its only realization.

Let E^{Rq} , q = 1, ..., Q be Q different but real variants of the collection of the introduced station signals. Let us also assume that the same number of numerical experiments, the homogenization adjustments, were performed and corresponding number of realizations of E^{H} were obtained using a chain of the calculations.

$$\boldsymbol{E}^{Rq} + \boldsymbol{X}^{C} = \boldsymbol{X}^{Iq}, \ \boldsymbol{X}^{Iq} \xrightarrow{g} \boldsymbol{X}^{Hq}, \ \boldsymbol{X}^{Hq} - \boldsymbol{X}^{C} = \boldsymbol{E}^{Hq}, \ q = 1, ..., Q$$
(6)

Based on these realizations, it is theoretically possible to evaluate $f_{M \times N}$. However, such task is hardly feasible in practice due to the extremely large number of dimensions to be considered. On the other hand, based on the



FIGURE 7 Examples of time series of errors: Real/introduced E_k^R (a), detected E_k^D (b) and residual E_k^H (c) calculated from the data presented in Figure 6

statistical ensemble of Q individual realizations of E^{H} we can evaluate some of the moments of the residual error distribution (5). In the context of our objective, the most important of them are a mean value (*m*) and some parameter that can characterize a width of the distribution such as the *SD* (σ) or the percentile range. The mean value provides information regarding the systematic bias of the homogenization adjustment, while the *SD* or the percentile range characterize its uncertainty. Both statistics, *m* and σ , can vary in the spacetime domain where E^{H} is defined and they can be evaluated using the following formulas

$$m_{ij} = \frac{1}{Q} \sum_{q=1}^{Q} e_{ij}^{Hq}, \qquad (7.1)$$

$$\sigma_{ij} = \left(\frac{1}{(Q-1)} \sum_{q=1}^{Q} \left(e_{ij}^{Hq} - m_{ij}\right)^2\right)^{\frac{1}{2}}, \quad (7.2)$$

While the proposed approach to the evaluation of the adjustment uncertainty on the daily time scale appears attractive and theoretically rigorous, it can potentially lead to some problems that may limit its practical applicability. For instance, one of the limitations can be related to difficulties with constructing a statistical ensemble for E^{R} with a sufficient number of its individual realizations in order to perform the calculations according to (6). Another example of possible limitations can be explained as follow: typically, at the end of the time domain T, all station signals in E^R contain undisturbed segments (see, e.g., Figure 7a). Hence, many of zero values in E^{H} are usually obtained there. Such zero values have to be excluded from the analysis when evaluating the homogenization adjustment since they do not mean the 'perfect' adjustment. However, it is not very easy to do so, because individual station signals usually have undisturbed segments of different length.

Estimating the statistical properties of the random field of the residual error E^{H} is not the only way to evaluate the

i = 1, ..., M, j = 1, ..., N.

performance of the homogenization adjustment and to quantify its uncertainty on the daily time resolution. An alternative approach is to use specially elaborated statistical metrics or indicators (e.g., Trewin, 2018; Vincent *et al.*, 2018). As noted in Coll *et al.* (2020), such metrics can provide useful indications in relation to the strengths and weaknesses of homogenization methods used.

2.3.2 | Metrics for the adjustment evaluation on the daily time scale

The performance evaluation of an adjustment algorithm and the quantification of its uncertainty are slightly different tasks in several aspects. For instance, we can evaluate the performance even if there is only a single realization of the adjustment output X^{H} . Whereas to define the uncertainty we normally should have a statistical ensemble of $X^{H}(X^{Hq}, q = 1, ..., Q)$ and the corresponding ensemble of $E^{H}(E^{Hq}, q = 1, ..., Q)$. As it was already mentioned above, a single realization of E^H can be used for the uncertainty quantification only if E^{H} satisfies special conditions. The evaluation is usually performed by means of some metrics or statistical indicators. The metrics are computed for each individual station in the data set based on error data E_i^H (i = 1, ..., M) or on comparison of the corresponding pair of time series X_i^H and X_i^C . Calculated for a single output of the homogenization adjustment X^{H} , they yield general (averaged in time) estimates of the systematic and random residual errors in this actual software run. The metrics values can be averaged over all stations, providing overall (for the whole space domain) evaluation. Some of such averaged metrics, however, can also be used to quantify the adjustment uncertainty.

Figure 8a shows a graphical comparison between the homogenized X_k^H and clean X_k^C time series, presented in Figure 6b,c. A similar plot for inhomogeneous X_k^I and clean X_k^C data (Figure 6a,c) is presented in Figure 8b for

comparison. The solid bisecting line of black colour, usually referred to as the line of true predictions, shows full agreement between respective time series. The perfect/ ideal adjustment algorithm would yield corrected values, which are exactly the same as the corresponding clean data. In this case, all dots depicting all pairs (x_{kj}^C, x_{kj}^H) , j = 1, ..., N would lie on the line of true predictions. The dots lying below the black line mean underestimation of the adjustment algorithm, while the dots above it show overestimation. Other lines in the diagrams are explained later. The figures are used below for further explanations.

The discrepancy between the homogenized and clean time series (Figure 8a) is obviously reduced compared to the discrepancy between the inhomogeneous and clean data (Figure 8b). The residual disagreement in Figure 8a might be quantified by means of some statistical metrics. Due to the random nature of X_k^H and X_k^C , it is evident, that several metrics should be used because no single one can provide complete information regarding the residual errors of both types, systematic and random.

Keeping in mind the daily resolution of our data, we applied six different metrics: the bias (B), root mean square error (RMSE), factor of exceedance (FOEX), percentage of days within $\pm 0.5/\pm 2^{\circ}$ C margin (POD05/ POD2), and difference in slopes (SlopeD). The use of metrics B, FOEX, and SlopeD is intended for estimating the systematic errors, while the other three, RMSE and POD05/POD2, are used for evaluation of the random or scatter residual errors. In the context of the uncertainty evaluation, the two most important metric are B and RMSE, which averaged values can also provide information regarding the overall deviation of the adjustment prediction from the true climate signal and the range of the possible residual errors, respectively. Formulas for most of the metrics are standard and well-known. However, we include them for completeness. Note that all formulas are presented for individual pairs of time series, X_i^H and X_i^C , i = 1, ..., M. Obviously, similar metrics can



FIGURE 8 Example of scatter diagrams. Homogenized \mathbf{x}_{K}^{H} (a) and raw \mathbf{x}_{K}^{I} (b) daily data are built against respective clean values \mathbf{x}_{K}^{C} presented in Figure 6

be calculated for inhomogeneous data by replacing X_i^H with X_i^I .

(1) Bias

$$B_{i} = \frac{1}{N_{i}} \sum_{j=1}^{N_{i}} \left(x_{ij}^{H} - x_{ij}^{C} \right) = \frac{1}{N_{i}} \sum_{j=1}^{N_{i}} e_{ij}^{H}, \qquad (8)$$

where N_i is a number of pairs (x_{ij}^C, x_{ij}^H) in an adjusted segment/segments. The data from the last uncorrected segment are not used in calculations $(N_i < N)$. The bias can be positive or negative. Depending on its sign it shows average overestimation (+) or underestimation (-) of the adjusted data. However, the bias does not provide any information whether overestimations are more frequent than underestimations or vice-versa. The 'perfect' homogenization algorithm would give 0 for this metric, while $B_i = 0$ does not mean that all differences $x_{ij}^H - x_{ij}^C = e_{ij}^H, j = 1, ..., N_i$ are zeros. In the case when a statistical ensemble of Q individual realizations of the adjustment outputs is available, B_i can be averaged over this statistical ensemble. By comparing (7.1) and (8) it becomes clear that such averaged value can be considered an estimate of the mean of the random field E^H for *i*-th station.

(2) RMSE

$$RMSE_{i} = \left(\frac{1}{N_{i}}\sum_{j=1}^{N_{i}} \left(x_{ij}^{H} - x_{ij}^{C}\right)^{2}\right)^{\frac{1}{2}} = \left(\frac{1}{N_{i}}\sum_{j=1}^{N_{i}} \left(e_{ij}^{H}\right)^{2}\right)^{\frac{1}{2}}.$$
 (9)

RMSE provides information about the average deviation of the adjusted data from the true climate signal. However, this metric can be also interpreted as a value that is proportional to the Euclidian distance between X_i^H and X_i^C in a multidimensional space. Consequently, such an interpretation provides qualitative explanation why *RMSE_i*, averaged over the statistical ensemble of *Q* model runs, can characterize the width of possible residual error distribution for the *i*-th station and, hence, can be used to characterize the homogenization adjustment uncertainty. Comparing (7.2) and (9), it can be concluded that such averaged value should be close to the *SD* of the random field E^H for the *i*-th station.

(3) Factor of exceedance

$$FOEX_{i} = \left(\frac{N_{(x_{ij}^{H} > x_{ij}^{C})}}{N_{i}} - 0.5\right) 100,$$
 (10)

where $N_{(x_{ij}^{H} > x_{ij}^{C})}$ is a number of pairs (x_{ij}^{C}, x_{ij}^{H}) when $x_{ij}^{H} > x_{ij}^{C}$, i.e., a homogenized value is overestimated in comparison with the respective value from a clean time series. The factor of exceedance is measured in percentage and its values range from -50 to 50%. For instance, FOEX = 50% means that all homogenized data are

overestimated with respect to the true climate data. This measure is widely used in climate analysis and applied meteorology, for example, Mosca *et al.* (1998).

(4, 5) Percentage of days within $\pm 0.5/\pm 2^{\circ}$ C margin. In addition to the line of true values in Figure 8, other reference lines might be shown on a scatter diagram in order to facilitate the qualitative evaluation of adjustment performance. For instance, pairs of parallels can be drawn that are defined as

$$\left|x^{H} - x^{C}\right| = \Delta T,\tag{11}$$

where |...| denotes an absolute value, x^H and x^C stand for independent and dependent variables (abscissas and ordinates) in Figure 8, respectively, ΔT is a certain threshold of temperature differences. Following Vincent *et al.* (2018), in our study as the thresholds we chose 0.5 and 2°C by analogy with the factor of 2 used in other fields of applied meteorology (e.g., Mosca *et al.*, 1998). A pair of such reference lines when $\Delta T = 2$ are shown in red colour in Figure 8. Now metrics *POD*05 and *POD*2 can be simply explained as percentage of dots $\left(x_{ij}^C, x_{ij}^H\right)$, which lie in the area between respective reference lines (11). That is,

$$POD05_{i} = \frac{N_{ij} |x_{ij}^{H} - x_{ij}^{C}| < 0.5}{N_{i}} 100 \text{ and } POD2_{i} = \frac{N_{ij} |x_{ij}^{H} - x_{ij}^{C}| < 2}{N_{i}} 100,$$
(12)

where $N_{|x_{ij}^H - x_{ij}^C| \le 0.5}$ and $N_{|x_{ij}^H - x_{ij}^C| \le 2}$ stand for the numbers of dots (x_{ij}^C, x_{ij}^H) , which lie in the areas inside respective lines (11). Such metrics characterize the magnitude of the scatter of the adjusted values around the clean data.

(6) Difference in slopes

$$SlopeD_i = b_i - 1, \tag{13}$$

where b_i is the slope of a linear regression model $X_i^H = a_i + b_i X_i^C$, which is built using the standard leastsquares approach. The need to introduce such metric can be explained based on Figure 8a. As can be seen from this figure, neither *B* nor *FOEX* can clearly capture the tendency of general simultaneous underestimation of positive temperatures and overestimation of negative ones (the opposite situation is also possible). The absolute values of the under/over-estimations depend on the temperature magnitude, and they are the largest for temperature extreme. In other words, the under/over-estimation should be reflected in the underestimation of the amplitude of the seasonal cycle showing less variability of the adjusted temperature values. We propose to evaluate such type of discrepancies (systematic error) between homogenized and clean data based on comparison of slopes of the true value line, which always equals to 1, and the linear regression built on the data (blue line in Figure 8). The metric is important when evaluating the adjustment of the daily data, since the under/overestimation of values from tails of the temperature distribution can affect the calculations of some climate extremes indices. The best value for SlopeD is 0. It is worth noting that a similar approach was used in (Della-Marta and Wanner, 2006), where a comparison of the candidate and reference series by means of a scatter diagram was part of the proposed adjustment method. According to that work, deviation of the slope of a line that fits the data from 1 indicates that daily temperatures at the candidate are less/more variable than those at the reference.

The set of the introduced metrics is capable of providing a fairly detailed description of the adjustment performance on the daily time resolution.

2.3.3 | Quantifying discrepancies between homogenized and clean data on the yearly scale

As it was pointed out in the introduction, daily air temperature data are used for computing climate extremes indices. Therefore, it is important to evaluate how the accuracy of the adjustment algorithm for data with such temporal resolution is reflected in calculation of these indices and their regular tendencies (trends) (Trewin and Trevitt, 1996). To do so, we calculated the yearly time series of the temperature data, TNy and TXy, and the following indices (Klein Tank et al., 2009; Zhang et al., 2011): FD (frost days), TR (tropical nights), TN10p (cold nights), TN90p (warm nights), ID (ice days), SU (summer days), TX10p (cold days), TX90p (warm days). However, due to peculiarities of the southern Sweden climate (relatively cold) we slightly shifted the standard absolute thresholds in the respective climate extremes indices. That is, instead of 0 and 20°C for FD and TR, respectively, we used -10 and 10°C. Instead of 0 and 25°C for ID and SU, respectively, the thresholds of 5 and 20°C were used. In order to indicate these changes in the calculating algorithms of the indices clearly, we will denote them as FD-10, TR10, ID5, and SU20. Calculation of the indices was performed for raw, clean, and homogenized data based on the RClimDex software (Zhang et al., 2018). After that, quantifying the discrepancies between the indices calculated based on the clean and

homogenized data was performed by means of only two metrics, namely *B* and *RMSE*. Similarly, to the daily time series, the metrics were calculated using the adjusted segment/segments only. In addition, we computed differences/errors in the indices linear trends (*TrD*), calculated for the adjusted and clean data. The trends were evaluated over the whole time series (including undisturbed segments) through the least squares regression.

2.3.4 | Ensemble of introduced station signals

As was noted above, the main source of the uncertainty for the homogenization adjustment is the station signals introduced into the raw time series. In other words, the results of the adjustment are sensitive to the input data and magnitude of errors contained there. It is natural to expect that the larger the deviation of raw time series from the clean ones, the larger the residual errors should be after the adjustment. In turn, the deviation of the raw time series from the clean data is controlled by the system of break points and corresponding statistical properties of homogeneous segments in the station signals E^R , such as the shift amplitudes/factors, signal to noise ratios etc. In real situation when homogenizing a some set of raw time series, such information is usually unknown. This is a reason why in order to estimate the adjustment uncertainty we have to use the benchmark data and consider all possible but real variants of the station signals or, in other words, consider their statistical ensemble E^{Rq} , q = 1, ..., Q.

Such ensemble is preferred for further calculations, no matter what approach is used to quantify the adjustment uncertainty: the statistical metrics or the random field formalism. Our general idea regarding creating E^{Rq} , q = 1, ..., Q is to use the collections of the error time series, introduced in the benchmark, and apply to them replacements and/or permutations. As was shown in Section 2.2., the collection of the station signals E^{R} , that was created in the INDE-CIS project, possesses statistical properties, which are close to reality. Therefore, we should expect that a large enough number of the replacements and/or permutations in the set of 94/96 (TN/TX, see Figure 3) different station signals will provide a sufficient number of individual realizations of E^{H} . Our methodology will be applied to two different case studies, with increasing complexity, which will be fully described in the Results section.

3 | RESULTS

3.1 | Case study #1

This first case study considers 10 stations (Figure 9) and limits the length of the corresponding time series to the period of 1971–1980 (10 years, similar to Vincent *et al.* (2018)). Nine time series (the references), belonging to the stations marked in black colour in Figure 9, are left clean, while the time series of the tenth station (the candidate), depicted in red, is assumed to be corrupted with only one break point dated to 01.01.1976. That is, the first half (1971–1975) of the period under study is intended to be corrupted. Using the same matrix notations as in (2), these initial conditions can be written as follows

$$\left\{x_{ij}^{I}\right\} = \left\{x_{ij}^{C}\right\}, \text{ when } i = 1, ..., 9, j = 1, ..., 3653,$$

or $i = 10, j = 1827, ..., 3653;$ (14.1)

$$\left\{x_{ij}^{I}\right\} \neq \left\{x_{ij}^{C}\right\}$$
, when $i=10, j=1, ..., 1826$, (14.2)

where 3653 is the total number of days in the time interval 1971–1980, while 1826 is the number of days in the interval 1971–1975.

The average distance between the candidate and reference stations is ~34 km, while the averaged Pearson's correlation coefficient between X_{10}^C and X_i^C , i = 1, ..., 9 is 0.96 for TN and 0.97 for TX data. Before the correlation calculation, the seasonal cycle was removed from each time series using an approach similar to Vincent *et al.* (2018).

In order to construct the raw data with the corrupted 5-year sub-period $(\left\{x_{ij}^{I}\right\}, i = 10, j = 1, ..., 1826)$, we analysed all station signals in \boldsymbol{E}^{R} , that were initially introduced in the INDECIS benchmark, and defined homogeneous error segments that have the length of more than 5 complete consecutive years (since January 1 until December 31). For instance, in the error time series shown in Figure 7 a, all three homogeneous non-zero segments, that is, (January 1, 1950-August 13, 1966), (August 14, 1966-February 19, 1972), (February 20, 1972-September 8, 2000), satisfy this condition. The total numbers of such segments in TN and TX error data sets are 185 and 193, respectively. Then 185 for TN and 193 for TX different versions of the raw time series were constructed by shifting (translating along the time axis) a 5-year period from each of the defined segments to 1971-1975 and adding them to the respective clean data $\{x_{ij}^{C}\}, i = 10, j = 1, ..., 1826$. This way (by performing such replacements), we obtained a statistical ensemble of individual realizations of the raw data set X^{Iq} , q = 1, ...,



FIGURE 9 The chosen set of meteorological stations in case study #1. Black dots show the stations whose time series were always clean, red square is the station where inhomogeneities were introduced

Q, where Q = 185 for TN and Q = 193 for TX. The members of the ensemble differ from each other by only statistical properties of the disturbed segment in the tenth series (see (14.1) and (14.2)), which are well known (Figures 4 and 5) and, hence, can be considered as controlled. Applying Climatol with the predefined break point to each member of the statistical ensemble, we obtained a sample of the respective number of the adjustment results, which were used for further calculations. It should be mentioned that the average correlation between X_{10}^{Iq} , q = 1, ..., Q and the system of the reference series X_i^C , i = 1, ..., 9 slightly varies for different q. For TN data the range of the correlation coefficient values is (0.80, 0.95) with the mean around 0.89, while for TX data the range and the mean are (0.81, 0.96) and 0.91, respectively. We believe that such variations are not substantially influencing on the adjustment results and, furthermore, they are unavoidable since they come from the variations of station signals in the statistical ensemble of the candidate time series.

The same corrupted period along with unchanged system of reference series allows to conduct statistically reliable and justified evaluation of the residual errors. Moreover, the approach, used in Case study #1, provides an assessment of a nearly pure effect of the introduced station signals on the adjustment uncertainty. This is because any other factors, which might have some effect on the homogenization adjustment, were kept approximately constant or removed.

Figure 10 shows the results of the adjustment uncertainty quantification on the daily scale by applying the concept of a random field to the residual errors E^{H} . Since only a single time series of the raw data set was corrupted on 1971–1975, E^{H} has non-zero values only for one point in the space domain (i.e., for tenth station) and just for the first half of the period under study. Therefore, the statistical properties of E^{H} were defined only for this station and period. In Figure 10, the mean values, 5th (P05) and 95th (P95) percentiles of empirical distributions of E^{H} , calculated for each day of 1971-1975, are shown. Figure (a) shows the calculations for TN, while (b) depicts the similar results for TX. The mean values were calculated based on formula (7.1), whereas the percentiles were evaluated using the samples of Q (185 for TN and 193 for TX) values e_{10i}^{Hq} , q = 1, ..., Q for each day (i = 1, ..., 1826).

As can be seen from the figure, the calculated parameters, means and percentiles, vary in time. Beside noise, which is due to the limited number of individual realizations in the statistical ensemble, a regular 1-year periodicity can be observed. Generally, the range of the residual error is less in summertime compared to winter months. Such non-stationary/periodic behaviour of the widths of the residual error distributions can be attributed to the similar periodicity of the introduced errors E^R . The reason for the seasonality in E^R is significantly less local spatial variability of air temperature in a summer period compared to winter. Thus, we could expect that the adjusted values of air temperatures, both TN and TX, are closer to the true climate signal in summer than in winter.

The similar 1-year periodicity of the mean values of the residual error distributions implies periodic bias of the air temperature, adjusted by the Climatol software. For both climatic variables, the residual errors are slightly shifted to negative values during summertime, while in winter months the shift has opposite direction. Such bias periodicity means the average underestimation of temperature in summer, and the overestimation in winter and it should have some influence on the amplitude of the seasonal cycle of the adjusted minimum and maximum air temperature.

In order to provide additional evidences for the conclusions, stated after the qualitative analysis of the results presented in Figure 10, we averaged the empirical



FIGURE 10 Mean, 5th and 95th percentiles (P05 and P95) of empirical distributions of the residual errors, evaluated for each day of the corrupted segment: (a) TN, (b) TX

error distributions over the whole period, and over January and July months separately (Figure 11). Table 1 contains some of the parameters of these averaged distributions. Similar parameters for the introduced arrors are presented in the table for comparison

duced errors are presented in the table for comparison. The seasonality of the residual error distributions is seen in the figure for both variables and it is also supported by the table content.

In summer months, the percentile intervals of the residual errors, (P05, P95), for the adjusted daily TN and TX air temperatures are (-2.80, 1.70) (°C) and (-2.60, 1.70)1.90) (°C), respectively. Note, that such quantitative assessments can be considered as one of possible measures of Climatol's adjustment uncertainty. The corresponding mean values of the error distributions are -0.41 and -0.22°C. These results imply that in summer we could expect any adjusted temperature value x_{ii}^{H} to be slightly underestimated (on average) compared to a respective clean temperature x_{ii}^C by 0.41°C for TN and 0.22°C for TX. In addition, we could expect with 90% probability that for minimum air temperature the x_{ii}^H adjusted value lays in the interval $(x_{ij}^C - 2.80, x_{ij}^C + 1.70)$ (°C), while for maximum air temperature the interval is $(x_{ij}^C - 2.60, x_{ij}^C + 1.90)$ (°C). It is important to note a reduction by ~26/11% (TN/TX) in the percentile range length of the residual errors compared to the introduced ones. Such decreasing of the uncertainty is a quantitative assessment of the added value (Sturm and Engström, 2019) of the homogenization International Journal <u>RMets</u> E2409

adjustment performed by the Climatol software on dayto-day level in a summer period.

In winter months, the ranges (P05, P95), evaluated for the homogenization adjustment errors in TN and TX data are (-3.60, 4.50) (°C) and (-2.00, 2.60) (°C), respectively. The corresponding mean values of the error distributions are 0.40°C for TN and 0.28°C for TX. Thus, in winter we could expect any adjusted temperature value x_{ii}^{H} to be slightly overestimated (on average) by 0.40°C for TN and 0.28°C for TX relatively to the respective clean value x_{ii}^C and with 90% probability it lays in the interval $(x_{ii}^C - 3.60, x_{ij}^C + 4.50)$ (°C) in case of TN air temperature and $\left(x_{ij}^{C}-2.00, x_{ij}^{C}+2.60\right)$ (°C) in case of TX. Compared with summer months, there is noticeable difference between widths of (P05, P95) intervals calculated for TN and TX winter residual errors. For minimum air temperature such interval is substantially larger (almost twice) meaning larger uncertainty in the adjusted values of TN in this period of the year. Similar to the summer period, the homogenization adjustment reduced the width of the introduced error distribution by 15/13% (TN/TX).

The parameters of the empirical distribution of the residual errors, averaged over the whole 5-year period (see Table 1), can characterize only overall (time-averaged) Climatol performance and uncertainty. Some peculiarities of the errors time evolution are neglected. For instance, the shifts of the error mean values in the opposite directions during the winter and summer seasons compensate each other yielding perfect, almost unbiased Climatol's adjustment. The 5th and 95th percentile for



FIGURE 11 Empirical distributions of the residual errors, averaged over (a, d) the whole 5-year period, (b, e) January months, (c, f) July months: (top panel) TN, (bottom panel) TX

		Year		January	January		July	
		$\overline{\boldsymbol{E}^{H}}$	E^R	$\overline{\boldsymbol{E}^{H}}$	E^R	$\overline{E^{H}}$	E^{R}	
TN	Mean	-0.03	-0.11	0.40	-0.08	-0.41	-0.13	
	SD	2.15	2.53	2.56	2.97	1.39	1.85	
	P05	-3.20	-4.00	-3.60	-4.90	-2.80	-3.20	
	P95	3.20	3.70	4.50	4.60	1.70	2.90	
	P95-P05	6.40	7.70	8.10	9.50	4.50	6.10	
TX	Mean	-0.02	-0.00	0.28	-0.03	-0.22	0.04	
	SD	1.64	1.84	1.58	1.78	1.48	1.67	
	P05	-2.50	-2.70	-2.00	-2.70	-2.60	-2.50	
	P95	2.30	2.60	2.60	2.60	1.90	2.50	
	P95-P05	4.80	5.30	4.60	5.30	4.50	5.00	

SKRYNYK ET AL.

TN and TX are between the respective summer and winter values, showing averaged uncertainty of the Climatol software. The *SD* of the residual error distributions, which also can be used to characterize the adjustment uncertainty along with the percentile range, are 2.15° C for TN and 1.64° C for TX. These numbers are important because they can be compared later with averaged values of *RMSE*, which are also intended to show the general/ overall uncertainty of the homogenization adjustment.

Thus, we can conclude, that if it is possible, the errors of the homogenization adjustment of daily air temperature time series should be evaluated on daily or, at least, seasonal scale. The overall time-averaged evaluation might omit some peculiarities of the residual errors. Figure 12 summaries evaluating results of Climatol's adjustment performance (including its uncertainty), which were obtained by applying the statistical metrics. It is important to keep in mind when interpreting these results that the metrics can provide only information regarding overall time-averaged performance of the software. As was pointed above, the six metrics that were used in the study yield detailed evaluation of Climatol's capability of removing systematic and random errors in each individual realization of the raw time series of a statistical ensemble. However, only averaged value of *RMSE* (averaged over a statistical ensemble) can be considered as a measure of the adjustment uncertainty, providing information regarding the width of empirical distribution



FIGURE 12 Boxplots of the metrics, calculated in the set of numerical experiments #1: (a) TN, (b) TX

of the potential residual errors. For each metric, 185/193 (TN/TX) values were calculated, that corresponds to the numbers of individual realizations in the statistical ensembles. These metric values are summarized as boxplots in the figure. Note, that the boxplots of the metrics, calculated for the respective raw data, are also shown for relative evaluation of the adjustment efficiency. Due to very short adjusted period (just 5 years) the climate extremes indices were not calculated and the evaluation of the Climatol software on the yearly scale was not performed in this series of numerical experiments.

As can be seen from the figure, the mean value of bias (*B*) and its interquartile range (IQR), which we use as a convenient measure of the metric distribution width directly shown in the boxplots, tend to zero for both variables, TN and TX. Similar tendencies are observed for *FOEX*. Here IQR is not zero, but it has relatively small magnitude, especially for TN. Both these metrics indicate the almost perfect performance of the Climatol software in removing systematic errors (shifts in the means). Such conclusion is plainly and brightly supported by a simple visual comparison with the same metrics in the raw data.

However, another type of the systematic residual errors associated with the seasonality of discrepancies between the homogenized and clean data (described by *SlopeD*) is not removed. Moreover, such type of errors seems to be slightly amplified by Climatol in a sense that almost all values of *SlopeD* became negative compared with the symmetric distribution of the metric values in the raw data. That means the simultaneous underestimation of summer temperatures and overestimation of winter ones, and as the result—the underestimation of the amplitude of seasonal cycle. Such conclusion is fully supported by the day-to-day evaluation provided above. The potential ability of the Climatol software to slightly alter seasonality was also pointed out by (Sturm and Engström, 2019).

The performance of the Climatol software in removing random errors is not so pronounced as the removing systematic ones. After adjusting, the means and IQRs of metrics *RMSE*, *POD*05, and *POD*2 for both variables, TN and TX, are slightly improved compared to similar values in the raw data. However, this improvement seems to be associated with the almost perfect removing of break point shifts in the means, and not directly related to the real Climatol's capability of coping with the scatter of errors. The mean value of *RMSE*, which yields the overall, time-averaged assessment of the adjustment uncertainty, is 2.06° C for TN and 1.53° C for TX. Such values are very close to the previously calculated *SD* of the residual error distributions, calculated on the day-to-day level and averaged over 5-year period (see Table 1). The coincidence of the uncertainty estimates that were obtained by applying different approaches indicates robustness of the drawn conclusions and the quantitative assessments. In addition, our assessments of *RMSE* for TN and TX adjusted data are close to similar estimates presented by Vincent *et al.* (2018).

It is worth noting again that the provided quantitative assessments of Climatol's performance and uncertainty (as well as those given in the following section) are valid only for cases when the correlation between candidate and reference series is quite high, ~ (0.80, 0.95) for TN and (0.81, 0.96) for Tx. As it was already mentioned, the uncertainty quantification in other situations, that is, with other values of correlation ties between time series, will be performed in our future work.

According to Vincent et al. (2018), adjustment algorithms, applied to daily air temperature data, might show worse capability of removing small size shifts compared to large ones. Thus, it would be interesting to define if there are some relationships between statistical characteristics of the introduced errors, such as their mean value (an amplitude of shift in the break point) and SD, and the corresponding values of the metrics, calculated after applying Climatol. The main purpose of the following calculations is to define what kind of errors (with small or large shift amplitude, with small or large noise component) is removed better. Because the statistical ensemble of Climatol runs for TN data contains 185 different individual realizations, the same number of different values of the error means and SD were calculated and bound to the corresponding values of the metrics (Figure 13). A similar figure was created also for TX, but it is not included in the text. Note, that in Figure 13 the metrics calculated based on the raw data are also shown for comparison.

The relationships for *B* and *FOEX* are trivial and they were expected due to the almost perfect performance of the Climatol software in removing jumps in the means. However, other metrics show more interesting dependencies on the error means and *SD*. For instance, *SlopeD* has negative values for any shift amplitude. However, the metric depends almost linearly on *SD* of the introduced errors. The larger the *SD*, the larger negative value of *SlopeD* should be expected, meaning the more intensive seasonality in the residual error time series. There are no any visible relations between the values of *RMSE*, *POD*05, and *POD2* and the shift amplitudes from some interval around zero (shifts of small magnitudes). In this interval (approximately from -2 to 2°C for TN and from



FIGURE 13 Relationships between the metric values and the main statistical properties of corrupted segment in the station signals: means (left column) and *SD* (right column). TN data. Red and blue colours mean homogenized/adjusted and raw data, respectively

-1 to 1°C for TX), there are also no visible differences between the metric values computed based on the homogenized and raw data. It means that removing shifts of small magnitudes has small influence on the random part of the residual errors. However, certain improvement of the metrics is observed for relatively large shifts. This conclusion is agreed well with the results by Vincent *et al.* (2018). Similar to *SlopeD*, the metrics *RMSE*, *POD*05, and *POD2* show noticeable relationships with the *SD* of the introduced errors. The larger magnitude of this statistical parameter, the larger random residual errors should be expected, what is indicated by the worse values of the metrics.

3.2 | Case study #2

This case study is more complex since the raw time series can have more than one break point and their positions are not strictly fixed: they are different in different realizations of the experiment. Here, we used the same 10 stations presented in Figure 9 but considered them on the initially defined period of time 1950–2005. Similar to Case study #1, nine time series (the references) are always kept clean, while constructing of the tenth disturbed or candidate series was slightly changed. Formally, these initial conditions can be stated in the following form

$$\left\{ x_{ij}^{I} \right\} = \left\{ x_{ij}^{C} \right\}, \text{ when } i = 1, ..., 9, j = 1, ..., 20454,$$

or $i = 10, j = N_{10} + 1, ..., 20454;$ (15.1)

$$\left\{x_{ij}^{I}\right\} \neq \left\{x_{ij}^{C}\right\}$$
, when $i=10, j=1,...,N_{10}$, (15.2)

where 20454 is the total number of days in the time interval 1950–2005, while N_{10} is the number of days in a disturbed segment/s of the candidate time series. N_{10} varies in different realizations of the numerical experiment.

In the INDECIS benchmark, 94 and 96 different nonzero station signals were created for TN and TX data, respectively (Figure 3). By adding these error series to the clean data of the tenth station alternately, we created corresponding numbers of different realizations of the raw data, which were used as inputs for the Climatol software. As in the previous case, each realization of this statistical ensemble consists of nine clean and one perturbed time series. By performing such replacement of the station signals, we do not change significantly the statistical properties of the introduced errors: the distributions of their means and SD are almost the same as in Case study #1. Besides, we do not change the system of reference stations. Pearson's correlation coefficients between X_{10}^C and X_i^C , i = 1, ..., 9 and between X_{10}^{Iq} (q = 1, ..., Q) and X_i^C , i = 1, ..., 9 are almost the same as in the previous case for both TN and TX data. But we change the structure and timing of break points (which positions are predefined during Climatol calculations), make it more difficult for the software to adjust different segments happened simultaneously in the raw time series. In addition, in this set of numerical experiments we can estimate Climatol's performance and its uncertainty on the yearly scale by defining the residual errors in the adjusted time series of the climate extremes indices. Evaluation of the Climatol software in case study #2 on the daily scale was performed only through metrics, that is, only overall, time-averaging evaluation was carried out. Day-to-day estimation of the residual error distributions, based on the concept of a random field, was not conducted. Such estimation is difficult to perform statistically correct in Case study #2 since individual realizations of the raw candidate time series in the statistical ensemble have last undisturbed periods of different lengths. Consequently, for days in the end of 1950-2005 calculations

would operate with considerably less quantity of nonzero error values compared with days in the beginning of 1950–2005.

Figure 14 contains boxplots of the metrics that were calculated on the daily scale for the adjusted TN and TX data. Similar to the previous case, we provided also corresponding metric values for raw data in order to evaluate relative success of the adjustment algorithm.

As it can be seen from the figure, the distributions of the metric values are almost the same as in the previous case. That means good Climatol's performance in removing systematic errors (shifts in the means) and moderate improvement of the metrics showing removing of scatter/ random residual errors. However, the seasonality of residual errors and the related issue of the underestimation of the seasonal cycle amplitude is also preserved in this case study. Therefore, a number of break points in the raw time series does not influence significantly the accuracy of Climatol's homogenization adjustment. If they are correctly defined during the detection process, the same (on average) adjustment results should be expected, no matter how many breaks were detected in each of raw time series.

The mean value of *RMSE* for the adjusted TN data is 2.07° C, while for the TX adjusted time series this parameter equals to 1.54° C. These values are very close to the similar estimates that were obtained in Case study #1. Thus, the overall time-averaged uncertainty of Climatol's adjustment is not influenced significantly by including multiple break points in the raw time series.

The boxplots of the metrics calculated based on the adjusted yearly time series of the air temperature data and the climate extremes indices are presented in Figure 15. Similar results that were obtained based on the raw yearly series are also presented in the figure for comparison. As can be seen in the figure, the averaging TN and TX daily data to the yearly scale almost completely remove the both types of residual errors. Nearly zero values of B for adjusted TNy and TXy series is a trivial result, since Climatol removes very well systematic errors even in daily data. The mean value of RMSE for TNy is reduced after adjustment from 0.94 to 0.20° C (by ~78%) while for TXy the reduction is slightly less: from 0.56 to 0.16°C (by ~63%). Such substantial improvement of RMSE for both climatic variables can be explained by the fact that averaging data to the yearly scale removes random/noisy part of the residual errors, seen on the daily scale. Note, that the mean values of RMSE, 0.20°C for TNy and 0.16°C for TXy, can be also considered as the measures of Climatol's adjustment uncertainty on the yearly time scale. In addition, as can be seen in the figure, Climatol removes most of the trend error in TNy and TXy data. The mean value and IQR of



FIGURE 14 Boxplots of the metrics calculated in the set of numerical experiments #2: (a) TN, (b) TX



FIGURE 15 Boxplots of the metrics calculated based on the yearly series of the climate extremes indices in the set of numerical experiments #2: (a) TN, (b) TX

TABLE 2 Parameters of empirical probability distributions of *TrD* (errors/differences in linear trends), defined for yearly time series of climate extreme indices: (a) TN, (b) TX

	FD-10, days/decade		TR10, days/decade		TN10p, %/decade		TN90p, %/decade	
(a)	Hom-cln	Raw-cln	Hom-cln	Raw-cln	Hom-cln	Raw-cln	Hom-cln	Raw-cln
Mean	0.29	-0.26	0.64	-0.79	-0.35	-0.52	-0.29	-0.73
SD	0.42	1.83	0.74	3.59	0.42	1.25	0.34	1.27
P05	-0.23	-3.00	-0.42	-6.65	-1.02	-2.22	-0.79	-2.54
P95	0.94	2.92	2.05	2.55	0.32	1.44	0.31	0.28
P95-P05	1.17	5.92	2.47	9.20	1.34	3.66	1.10	2.82
	ID5, days/decade		SU20, days/decade		TX10p, %/decade		TX90p, %/decade	
(b)	Hom-cln	Raw-cln	Hom-cln	Raw-cln	Hom-cln	Raw-cln	Hom-cln	Raw-cln
Mean	0.05							
	-0.05	-0.36	0.21	-0.56	-0.13	-0.13	-0.10	-0.36
SD	-0.05 0.27	-0.36 0.88	0.21 0.44	-0.56 1.73	-0.13 0.33	-0.13 0.79	-0.10 0.23	-0.36 0.64
SD P05	-0.05 0.27 -0.49	-0.36 0.88 -1.88	0.21 0.44 -0.37	-0.56 1.73 -3.41	-0.13 0.33 -0.71	-0.13 0.79 -1.47	-0.10 0.23 -0.49	-0.36 0.64 -1.40
SD P05 P95	-0.05 0.27 -0.49 0.39	-0.36 0.88 -1.88 0.96	0.21 0.44 -0.37 0.96	-0.56 1.73 -3.41 2.00	-0.13 0.33 -0.71 0.33	-0.13 0.79 -1.47 1.06	-0.10 0.23 -0.49 0.23	-0.36 0.64 -1.40 0.56

TrD are almost zeros (~0.00 and ~ 0.01° C/decade, respectively) for both climatic variables.

Climatol removes well both types of errors also in the time series of all considered extreme indices. This is clearly seen in the figure, where empirical distributions of *B* and *RMSE*, calculated based on the adjusted data, can be compared with similar distributions, obtained for the raw series. Both metrics for all indices indicate substantial improvement after applying Climatol's adjustment. The underestimation of the seasonal cycle amplitude in the adjusted data, seen on the daily time resolution, is not so noticeable in the indices time series, probably due to relatively small negative values of *SlopeD* (see Figure 14). However, the means of *B* for all indices with fixed thresholds are slightly negative, meaning general slight underestimation of these indices in the adjusted data.

Below we focus mainly on trend evaluation in the time series of the extreme indices due to their critical importance in climatological applications. The empirical distributions of errors (differences) in trends, *TrD*, calculated for the adjusted data are also presented in Figure 15. Table 2 contains some of parameters of the empirical distributions of *TrD* values. The first noticeable qualitative conclusion that can be drawn from the figure is substantial decreasing of the trend errors in the adjusted data compared with the raw ones. Regular tendencies of all extreme indices, evaluated based on the corrected data, are much closer to real trends than evaluated based on the raw time series.

Based on the table content, quantitative assessments of Climatol's accuracy and uncertainty in the indices trend calculation can be derived. For instance, the mean value of the trend errors in the adjusted series of FD-10 (frost days) is relatively small, 0.29 days/decade (2.9 days/100years). The uncertainty of the trend calculation in the adjusted FD-10 data can be estimated by mean of the SD (0.42 days/decade) or the percentile range (P05, P95), which is (-0.23, 0.94)(days/decade). Thus, we could expect, that a linear trend, calculated in the FD-10 yearly time series that was corrected by the Climatol software, is slightly shifted (on average) on 0.29 days/decade relatively to the true climate trend (Tr^{C}) , and with 90% probability it lie in the interval $(Tr^{C} - 0.23, Tr^{C} + 0.94)$ (days/decade). It is worth noting, that the percentile range of the trend errors in the raw time series is significantly larger, (-3.00,2.92)(days/decade), that is, after applying Climatol, a 80% decrease of the uncertainty can be reported. Similar assessments can be obtained from Table 2 for other climate extreme indices. We also can conclude, that, in general, trends can be estimated more accurately and with less uncertainty in the adjusted time series of the TX extreme climate indices than in TN extremes. One more important conclusion is that despite the substantial amount of the residual scatter/random errors which still remain in the adjusted daily time series, the linear trends calculated on the corrected yearly time series are reliable and close to real regular tendencies and they can be evaluated with significantly removed uncertainty.

4 | CONCLUSION

In this study, the uncertainty quantification and the general performance evaluation of Climatol's adjustment algorithm, applied to daily minimum and maximum air temperature time series, are presented. We focused our attention only on the most influencing and important source of the uncertainty, namely introduced station signals into the raw data set to be adjusted. Other possible sources of the adjustment uncertainty were removed from the analysis or kept approximately constant. For instance, the mean correlation between candidate and reference series was around (0.80, 0.95) for TN and (0.81, 0.96) for Tx data. Therefore, our results are valid only for cases where the mentioned mean correlation can be observed. The sensitivity of the obtained quantitative assessments to other factors/sources will be addressed in our future work.

In order to evaluate the adjustment uncertainty, we used the INDECIS benchmark data and applied a complex approach, quantifying the uncertainty at different levels of detail and time resolution. According to our findings, Climatol's adjustment uncertainty, evaluated on day-to-day level, varies in time, and depends on the season. In summer months, the residual errors in the adjusted daily TN and TX series are expected to belong to the intervals, (P05, P95), (-2.80, 1.70) (°C) and (-2.60, 1.90) (°C), respectively. In winter months, the ranges of the possible remaining errors are larger: (-3.60, 4.50)(°C) for TN and (-2.00, 2.60) (°C) for TX. The overall adjustment uncertainty, averaged over all seasons, can be evaluated as the error range, (P05, P95), (-3.20, 3.20)(°C) for TN and (-2.50, 2.30) (°C) for TX. In terms of SD of the residual error distributions, the overall uncertainty can be evaluated as 2.15°C for TN and 1.64°C for TX data. These estimates agree well with the mean values of, which also can be used as a measure of the width of the empirical distribution of the residual errors. Besides 1-year periodicity in the width of the residual error distributions, their mean values are also slightly shifted periodically. For both climatic variables, the shift is towards negative values during summertime, while in winter months it has opposite direction. Such peculiarities of the residual errors can lead to the slight underestimation of the amplitude of the seasonal cycle of the adjusted TN and TX data. The calculations based on the specially introduced metric (SlopeD) provide additional evidence for such conclusion. Other metrics, used in the study, showed that Climatol removes extremely well systematic errors related to jumps in the mean and this Climatol's capability is valid for shifts of any magnitude and does not depend on the number of break points in the raw time series. The ability of Climatol to remove scatter/random errors in the daily raw time series is not so pronounced.

However, on the yearly time scale, both types of residual errors are removed well in adjusted time

series. The adjusted yearly TN and TX temperature data are unbiased, and their uncertainty is reduced significantly: mean values of *RMSE* for TNy and TXy were decreased to 0.20° C (by ~78%) and 0.16° C (by ~63%), respectively. In addition, Climatol removes most of the trend error in TNy and TXy data, so trend analysis is more solid and better represents climate variations.

Similar conclusions are valid for the yearly time series of the considered climate extreme indices: both types of errors are removed well by Climatol. The underestimation of the seasonal cycle amplitude in the adjusted data, seen on the daily time resolution, is not clearly reflected in the indices time series. However, the mean values of bias (B) for all indices with fixed thresholds are slightly negative, meaning slight underestimation of these indices in the adjusted data. However, this does not have substantial influence on the linear trend calculations in the indices time series. The trends calculated in the adjusted time series are generally unbiased. The percentile (P05, P95) ranges of the errors in the indices trends, calculated based on adjusted data, is reduced by ~70-80% compared to the trend errors in the corresponding raw time series. Despite the substantial amount of the residual scatter errors in daily time series, the linear trends calculated on the corrected yearly time series are close to real regular tendencies and they can be evaluated with significantly removed uncertainty.

The next step to be undertaken in the context of Climatol's uncertainty evaluation is the analysis of how the quantitative assessments obtained depend on the correlation between a candidate and reference series. In addition, a similar assessment of Climatol's adjustment algorithm applied to daily precipitation data should be performed. Also, a much more difficult case, when the uncertainty of the adjustment and detection algorithms are evaluated simultaneously, should be considered in order to obtain the complete picture of Climatol capability to cope with inhomogeneities on the daily time scale.

ACKNOWLEDGEMENTS

The work was performed in the frame of the INDECIS project, that is a part of ERA4CS, an ERA-NET initiated by JPI Climate, and funded by FORMAS (SE), DLR (DE), BMWFW (AT), IFD (DK), MINECO (ES), ANR (FR) with co-funding by the European Union (Grant 690462). The work has been partially supported by the Ministry of Education and Science of Kazakhstan (Grant BR05236454) Nazarbayev University and (Grant 090118FD5345). The authors are grateful to anonymous reviewers for careful reading of the manuscript and valuable comments and suggestions they have made.

ORCID

Oleg Skrynyk https://orcid.org/0000-0001-8827-0280 Enric Aguilar https://orcid.org/0000-0002-8384-377X Jose Guijarro https://orcid.org/0000-0002-9527-9758 Luc Yannick Andreas Randriamarolaza https://orcid. org/0000-0002-2939-2250

Sergiy Bubin ^b https://orcid.org/0000-0002-2783-078X

REFERENCES

- Aguilar, E., Auer, I., Brunet, M., Peterson, T.C. and Wieringa, J. (2003) WMO Guidelines on Climate Metadata and Homogenization. WCDMP No.53, WMO-TD No. 1186. Geneva, Switzerland: WMO.
- Aguilar, E., van der Schrier, G., Guijarro, J.A., Stepanek, P., Zahradnicek, P., Sigro, J., Coscarelli, R., Engstrom, E., Curley, M., Caloiero, T., Lledo, L., Ramon, J. and Antonia Valente, M. (2018) *Quality Control and Homogenization Benchmarking-Based Progress from the INDECIS Project*. Vienna, Austria: General Assembly of the European Geosciences Union 8–13 April 2018, EGU2018-16392.
- Alexander, L.V., Zhang, X., Peterson, T.C., Caesar, J., Gleason, B., Klein Tank, A.M.G., Haylock, M., Collins, D., Trewin, B., Rahim, F., Tagipour, A., Kumar Kolli, R., Revadekar, J.V., Griffiths, G., Vincent, L., Stephenson, D.B., Burn, J., Aguilar, E., Brunet, M., Taylor, M., New, M., Zhai, P., Rusticucci, M. and Vazquez Aguirre, J.L. (2006) Global observed changes in daily climate extremes of temperature and precipitation. *Journal of Geophysical Research*, 111, D05109. https://doi.org/10.1029/2005JD006290.
- Alexandersson, H. (1986) A homogeneity test applied to precipitation data. Journal of Climatology, 6(6), 661–675. https://doi.org/ 10.1002/joc.3370060607.
- Alexandersson, H. and Moberg, A. (1997) Homogenization of Swedish temperature data. Part I: homogeneity test for linear trends. *International Journal of Climatology*, 17(1), 25–34. https://doi. org/10.1002/(SICI)1097-0088(199701)17:1<25::AID-JOC103>3. 0.CO;2-J.
- Azorin-Molina, C., Guijarro, J.A., McVicar, T.R., Trewin, B.C., Frost, A.J. and Chen, D. (2019) An approach to homogenize daily peak wind gusts: an application to the Australian series. *International Journal of Climatology*, 39(4), 2260–2277. https:// doi.org/10.1002/joc.5949.
- Brunet, M., Saladié, O., Jones, P., Sigró, J., Aguilar, E., Moberg, A., Lister, D., Walther, A. and Almarza, C. (2008) A Case-Study/Guidance on the Development of Long-Term Daily Adjusted Temperature Datasets. WMO-TD no. 1425, WCDMP no. 66. Geneva: World Meteorological Organization.
- Coll, J., Domonkos, P., Guijarro, J., Curley, M., Rustemeier, E., Aguilar, E., Walsh, S. and Sweeney, J. (2020) Application of homogenization methods for Ireland's monthly precipitation records: comparison of break detection results. *International Journal of Climatology*, 1–20. https://doi.org/10.1002/joc. 6575.
- Collins, W.J., Bellouin, N., Doutriaux-Boucher, M., Gedney, N., Hinton, T., Jones, C.D., Liddicoat, S., Martin, G., O'Connor, F., Rae, J., Senior, C., Totterdell, I., Woodward, S., Reichler, T. and Kim, J. (2008) Evaluation of the HadGEM2 model. *MetOffice Hadley Centre Technical Note*, 74, 47.

- DeGaetano, A.T. (2006) Attributes of several methods for detecting discontinuities in mean temperature series. *Journal of Climate*, 19(5), 838–853. https://doi.org/10.1175/JCLI3662.1.
- Della-Marta, P. and Wanner, H. (2006) A method of homogenizing the extremes and mean daily temperature measurements. *Journal of Climate*, 19(17), 4179–4197. https://doi.org/10.1175/ JCLI3855.1.
- Domonkos, P. (2011) Efficiency evaluation for detecting inhomogeneities by objective homogenization methods. *Theoretical and Applied Climatology*, 105, 455–467. https://doi.org/10.1007/ s00704-011-0399-7.
- Domonkos, P. (2017) Time series homogenization with optimal segmentation and ANOVA correction: past, present and future. Proceeding of 9th Seminar for homogenization and quality control in climatological databases and 4th conference on spatial interpolation techniques in climatology and meteorology (Budapest, April 3-7), WMO WCDMP-no.85, pp. 46-62.
- Domonkos, P. and Coll, J. (2017) Time series homogenization of large observational datasets: impact of the number of partner series on efficiency. *Climate Research*, 74, 31–42. https://doi. org/10.3354/cr01488.
- Domonkos, P. and Efthymiadis, D. (2013) Development and testing of homogenization methods: moving parameter experiments with ACMANT. *Advances in Science and Research*, 10, 43–50. https://doi.org/10.5194/asr-10-43-2013.
- Ducré-Robitaille, J.-F., Vincent, L.A. and Boulet, G. (2003) Comparison of techniques for detection of discontinuities in temperature series. *International Journal of Climatology*, 23(9), 1087–1101. https://doi.org/10.1002/joc.924.
- Dumitrescu, A., Cheval, S. and Guijarro, J.A. (2020) Homogenization of a combined hourly air temperature dataset over Romania. *International Journal of Climatology*, 40(5), 2599–2608. https://doi.org/10.1002/joc.6353.
- Fioravanti, G., Piervitali, E. and Desiato, F. (2019) A new homogenized daily data set for temperature variability assessment in Italy. *International Journal of Climatology*, 39(15), 5635–5654. https://doi.org/10.1002/joc.6177.
- Guijarro, J.A. (2011) Influence of network density on homogenization performance. Proceeding of 7th Seminar for Homogenization and Quality Control in Climatological Databases Jointly Organized with the Meeting of COST ES0601 (HOME) Action MC Meeting. Budapest, Hungary, 24-27 October, WMO WCDMP-No. 78, pp. 11-18.
- Guijarro, J.A. (2018) Homogenization of climatic series with Climatol. Version 3.1.1. Guide.
- Guijarro, J.A., Aguilar, E., Caloiero, T., Coscarelli, R., Curley, M. and Pérez-Zanón, N. (2018) Homogenization of daily Essential Climatic Variables with Climatol 3.1 within the INDECIS project. Budapest, Hungary: European Conference for Applied Meteorology and Climatology, 3-7 September 2018, EMS2018-413.
- Guijarro, J.A., Aguilar, E., Domoncos, P., Sigró, J., Štepánek, P., Venema, V. and Zahradnícek, P. (2019) Benchmarking Results of the Homogenization of Daily Essential Climatic Variables within the INDECIS Project. Vienna, Austria: General Assembly of the European Geosciences Union, pp. 7–12 April 2019, EGU2019-10896-1.
- Guijarro, J.A., López, J.A., Aguilar, E., Domonkos, P., Venema, V. K.C., Sigró, J. and Brunet, M. (2017) Comparison of

homogenization packages applied to monthly series of temperature and precipitation: The MULTITEST project. *Proceeding of* 9th Seminar for homogenization and quality control in climatological databases and 4th conference on spatial interpolation techniques in climatology and meteorology. Budapest, Hungary, 3-7 April 2017, WMO WCDMP-no.85, pp. 46-62.

- Hartmann, D.L., Klein Tank, A.M.G., Rusticucci, M., Alexander, L.
 V., Brönnimann, S., Charabi, Y., Dentener, F.J., Dlugokencky, E.J., Easterling, D.R., Kaplan, A., Soden, B.J., Thorne, P.W., Wild, M. and Zhai, P.M. (2013) Observations: atmosphere and surface. In: *Climate Change: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change.* Cambridge, UK and New York. NY: Cambridge University Press.
- Iman, R.L. and Helton, J.C. (1988) An investigation of uncertainty and sensitivity analysis techniques for computer models. *Risk Analysis*, 8(1), 71–90. https://doi.org/10.1111/j.1539-6924.1988.tb01155.x.
- INDECIS (2018) Integrated approach for the development across Europe of user oriented climate indicators for GFCS highpriority sectors: agriculture, disaster risk reduction, energy, health, water and tourism. Available at: http://www.indecis.eu/ [Accessed August 3, 2020].
- Jakeman, A.J., Letcher, R.A. and Norton, J.P. (2006) Ten iterative steps in development and evaluation of environmental models. *Environmental Modelling and Software*, 21(5), 602–614. https:// doi.org/10.1016/j.envsoft.2006.01.004.
- Killick, R.E. (2016) Benchmarking the performance of homogenization algorithms on daily temperature data. *PhD Thesis*, University of Exeter, 249 pp. Available at: https://ore.exeter.ac.uk/ repository/handle/10871/23095 [Accessed August 3, 2020].
- Klein Tank, A.M.G., Wijngaard, J.B., Können, G.P., Böhm, R., Demarée, G., Gocheva, A., Mileta, M., Pashiardis, S., Hejkrlik, L., Kern-Hansen, C., Heino, R., Bessemoulin, P., Müller-Westermeier, G., Tzanakou, M., Szalai, S., Pálsdóttir, T., Fitzgerald, D., Rubin, S., Capaldo, M., Maugeri, M., Leitass, A., Bukantis, A., Aberfeld, R., van Engelen, A.F.V., Forland, E., Mietus, M., Coelho, F., Mares, C., Razuvaev, V., Nieplova, E., Cegnar, T., Antonio López, J., Dahlström, B., Moberg, A., Kirchhofer, W., Ceylan, A., Pachaliuk, O., Alexander, L.V. and Petrovic, P. (2002) Daily dataset of 20th-century surface air temperature and precipitation series for the European climate assessment. *International Journal of Climatology*, 22(12), 1441–1453. https://doi.org/10.1002/joc.773.
- Klein Tank, A.M.G., Zwiers, F.W. and Zhang, X. (2009) Guidelines on analysis of extremes in a changing climate in support of informed decisions for adaptation, climate data and monitoring WCDMP-no 72, WMO-TD No 1500, p 55.
- Kuglitsch, F.G., Auchmann, R., Bleisch, R., Broennigmann, S., Martius, O. and Stewart, M. (2012) Break detection of annual Swiss temperature series. *Journal of Geophysical Research*, 117, D13105. https://doi.org/10.1029/2012JD017729.
- Lindau, R. and Venema, V. (2016) The uncertainty of break positions detected by homogenization algorithms in climate records. *International Journal of Climatology*, 36(2), 576–589. https://doi.org/10.1002/joc.4366.
- Mamara, A., Argiriou, A.A. and Anadranistakis, M. (2013) Homogenization of mean monthly temperature time series of Greece. *International Journal of Climatology*, 33(12), 2649–2666. https://doi.org/10.1002/joc.3614.

- Mamara, A., Argiriou, A.A. and Anadranistakis, M. (2014) Detection and correction of inhomogeneities in Greek climate temperature series. *International Journal of Climatology*, 34(10), 3024–3043. https://doi.org/10.1002/joc.3888.
- Meseguer-Ruiz, O., Ponce-Philimon, P.I., Quispe-Jofré, A.S., Guijarro, J.A. and Sarricolea, P. (2018) Spatial behavior of daily observed extreme temperatures in Northern Chile (1966–2015): data quality, warming trends, and its orographic and latitudinal effects. *Stochastic Environmental Research and Risk Assessment*, 32, 3503–3523. https://doi.org/10.1007/s00477-018-1557-6.
- Mestre, O., Gruber, C., Prieur, C., Caussinus, H. and Jourdain, S. (2011) SPLIDHOM: a method for homogenization of daily temperature observations. *Journal of Applied Meteorology* and Climatology, 50(11), 2343–2358. https://doi.org/10.1175/ 2011JAMC2641.1.
- Mosca, S., Graziani, G., Klug, W., Bellasio, R. and Bianconi, R. (1998) A statistical methodology for the evaluation of longrange dispersion models: an application to the ETEX exercise. *Atmospheric Environment*, 32(24), 4307–4324. https://doi.org/ 10.1016/S1352-2310(98)00179-4.
- Osadchyi, V., Skrynyk, O.A., Radchenko, R. and Skrynyk, O.Y. (2018) Homogenization of Ukrainian air temperature time series. *International Journal of Climatology*, 38(1), 497–505. https://doi.org/10.1002/joc.5191.
- Pérez-Zanón, N., Sigró, J., Aguilar, E., Guijarro J.A., van der Schrier, G., Stepanek, P., Zahradnicek, P., Coscarelli, R., Engström, E., Curley, M., Caloiero, T., Lledó, L., Ramon, J., Valente, M.A. and Carvalho, S. (2018) First steps towards a benchmarking experiment in quality control and homogenization of observed data. Budapest, Hungary: European Conference for Applied Meteorology and Climatology, 3-7 September 2018, EMS2018-465.
- Prohom, M., Barriendosb, M. and Sanchez-Lorenzod, A. (2016) Reconstruction and homogenization of the longest instrumental precipitation series in the Iberian Peninsula (Barcelona, 1786–2014). *International Journal of Climatology*, 36(8), 3072–3087. https://doi.org/10.1002/joc.4537.
- Reeves, J., Chen, J., Wang, X.L., Lund, R. and Lu, Q. (2007) A review and comparison of change points detection techniques for climate data. *Journal of Applied Meteorology and Climatol*ogy, 46, 900–915. https://doi.org/10.1175/JAM2493.1.
- Sanchez-Lorenzo, A., Wild, M., Brunetti, M., Guijarro, J.A., Hakuba, M.Z., Calbó, J., Mystakidis, S. and Bartok, B. (2015) Reassessment and update of long-term trends in downward surface shortwave radiation over Europe (1939–2012). *Journal of Geophysical Research-Atmospheres*, 120(18), 9555–9569. https:// doi.org/10.1002/2015JD023321.
- Skrynyk, O.Y., Aguilar, E., Skrynyk, O.A., Sidenko, V., Boichuk, D. and Osadchyi, V. (2019) Quality control and homogenization of monthly extreme air temperature of Ukraine. *International Journal of Climatology*, 39(4), 2071–2079. https://doi.org/10. 1002/joc.5934.
- Sokal, R.R. and Rohlf, P.J. (1969) *Introduction to Biostatistics*, 2nd edition. New York: W.H. Freeman, p. 363.
- Squintu, A.A., van der Schrier, G., Štěpánek, P., Zahradníček, P., Klein Tank, A. (2020) Comparison of homogenization methods for daily temperature series against an observation-based benchmark dataset. *Theoretical and Applied Climatology*, 140 (1-2), 285–301. http://dx.doi.org/10.1007/s00704-019-03018-0.

- Sturm, C. and Engström, E. (2019) Estimating the sensitivity and accuracy of homogenization: a case study with Climatol on temperature from the INDECIS benchmark. 12th EUMETNET Data Management Workshop, De Bilt, the Netherlands, 6–8 November 2019.
- Szentimrey, T. (2008) Methodological questions of series comparison. Proceeding of 6th Seminar for Homogenization and Quality Control in Climatological Databases. Budapest, Hungary, 26-30 May 2008, WMO WCDMP-No. 76, pp. 1–7.
- Trewin, B. (2010) Exposure, instrumentation, and observing practice effects on land temperature measurements. WIREs Climate Change, 1(4), 490–506. https://doi.org/10.1002/wcc.46.
- Trewin, B. (2013) A daily homogenized temperature data set for Australia. International Journal of Climatology, 33(6), 1510–1529. https://doi.org/10.1002/joc.3530.
- Trewin, B. (2018) The Australian Climate Observations Reference Network – Surface Air Temperature (ACORN-SAT).Version
 2. Bureau Research Report No. 032. Available at: http://www. bom.gov.au/climate/change/acorn-sat/documents/BRR-032. pdf. [Accessed August 3, 2020].
- Trewin, B.C. and Trevitt, A.C.F. (1996) The development of composite temperature records. *International Journal of Climatology*, 16(11), 1227–1242. https://doi.org/10.1002/(SICI)1097-0088 (199611)16:11<1227::AID-JOC82>3.0.CO;2-P.
- Venema, V., Mestre, O., Aguilar, E., Auer, I., Guijarro, J.A., Domonkos, P., Vertacnik, G., Szentimrey, T., Stepanek, P., Zahradnicek, P., Viarre, J., Muller-Westermeier, G., Lakatos, M., Williams, C.N., Menne, M., 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. *Climate of the Past*, 8, 89–115. https://doi.org/ 10.5194/cp-8-89-2012.
- Vincent, L.A., Milewska, E.J., Wang, X.L. and Hartwell, M.M. (2018) Uncertainty in homogenized daily temperatures and derived indices of extremes illustrated using parallel observations in Canada. *International Journal of Climatology*, 38(2), 692–707. https://doi.org/10.1002/joc.5203.

- Walker, W.E., Harremoës, P., Rotmans, J., van der Sluijs, J.P., van Asselt, M.B.A., Janssen, P. and Krayer von Krauss, M.P. (2003)
 Defining uncertainty: a conceptual basis for uncertainty management in model-based decision support. *Integrated Assessment*, 4(1), 5–17. https://doi.org/10.1076/iaij.4.1.5.16466.
- Willett, K., Williams, C., Jolliffe, I.T., Lund, R., Alexander, L.V., Brönnimann, S., Vincent, L.A., Easterbrook, S., Venema, V.K.C., Berry, D., Warren, R.E., Lopardo, G., Auchmann, R., Aguilar, E., Menne, M.J., Gallagher, C., Hausfather, Z., Thorarinsdottir, T. and Thorne, P.W. (2014) A framework for benchmarking of homogenization algorithm performance on the global scale. *Geoscientific Instrumentation, Methods and Data Systems*, 3, 187–200. https://doi.org/10.5194/gi-3-187-2014.
- Yosef, Y., Aguilar, E. and Alpert, P. (2018) Detecting and adjusting artificial biases of long-term temperature records in Israel. *International Journal of Climatology*, 38(8), 3273–3289. https:// doi.org/10.1002/joc.5500.
- Yozgatligil, C. and Yazici, C. (2016) Comparison of homogeneity tests for temperature using a simulation study. *International Journal* of *Climatology*, 36(1), 62–81. https://doi.org/10.1002/joc.4329.
- Zhang, X., Alexander, L., Hegerl, G.C., Jones, P., Klein Tank, A., Peterson, T.C., Trewin, B. and Zwiers, F.W. (2011) Indices for monitoring changes in extremes based on daily temperature and precipitation data. WIREs Climate Change, 2(6), 851–870. https://doi.org/10.1002/wcc.147.
- Zhang, X., Feng, Y., Chan, R. (2018) Introduction to RClimDex v1.9. Guide. Climate research Division, Environment Canada, Downsview Ontario, Canada.

How to cite this article: Skrynyk O, Aguilar E, Guijarro J, Randriamarolaza LYA, Bubin S. Uncertainty evaluation of Climatol's adjustment algorithm applied to daily air temperature time series. *Int J Climatol*. 2021;41 (Suppl. 1): E2395–E2419. <u>https://doi.org/10.1002/joc.6854</u>