



Impact of bias nonstationarity on the performance of uni- and multivariate bias-adjusting methods

Jorn Van de Velde^{1,2}, Matthias Demuzere^{3,1}, Bernard De Baets², and Niko E. C. Verhoest¹

¹Hydro-Climatic Extremes Lab, Ghent University, Ghent, Belgium

²KERMIT, Department of Data Analysis and Mathematical Modelling, Ghent University, Ghent, Belgium

³Department of Geography, Ruhr-University Bochum, Bochum, Germany

Correspondence: Jorn Van de Velde (jorn.vandavelde@ugent.be)

Abstract.

Climate change is one of the biggest challenges currently faced by society, with an impact on many systems, such as the hydrological cycle. To locally assess this impact, Regional Climate Model (RCM) simulations are often used as input for hydrological rainfall-runoff models. However, RCM results are still biased with respect to the observations. Many methods have been developed to adjust these biases, but only during the last few years, methods to adjust biases that account for the correlation between the variables have been proposed. This correlation adjustment is especially important for compound event impact analysis. As a simple example of those compound events, hydrological impact assessment is used here, as hydrological models often need multiple locally unbiased input variables to ensure an unbiased output. However, it has been suggested that multivariate bias-adjusting methods may perform poorly under climate change conditions because of bias nonstationarity. In this study, two univariate and three multivariate bias-adjusting methods are compared with respect to their performance under climate change conditions. To this end, the methods are calibrated in the late 20th century (1970-1989) and validated in the early 21st century (1998-2017), in which the effect of climate change is already visible. The variables adjusted are precipitation, evaporation and temperature, of which the former two are used as input for a rainfall-runoff model, to allow for the validation of the methods on discharge. Although not used for discharge modelling, temperature is a commonly-adjusted variable in both uni- and multivariate settings and therefore important to take into account. The methods are also evaluated using indices based on the adjusted variables, the temporal structure, and the multivariate correlation. For precipitation, all methods decrease the bias in a comparable manner. However, for many other indices the results differ considerably between the bias-adjusting methods. The multivariate methods often perform worse than the univariate methods, a result that is especially notable for temperature and evaporation. As these variables have already changed the most under climate change conditions, this reinforces the opinion that the multivariate bias-adjusting methods are not yet fit to cope with nonstationary climate conditions. Although the effect is slightly dampened by the hydrological model, our analysis still reveals that, to date, the simpler univariate bias-adjusting methods are preferred for assessing climate change impact.



1 Introduction

25 The influence of climate change is felt throughout many regions of the world, as becomes evident from the higher frequency or intensity of natural hazards, such as floods, droughts, heatwaves and forest fires (IPCC, 2012). As these intensified natural hazards threaten society, it is essential to be prepared for them. Knowledge on future climate change is obtained by running Global Climate Models (GCMs), creating large ensemble outputs such as in the Climate Model Intercomparison Project 6 (CMIP6) (Eyring et al., 2016). Although they are informative on a global scale, the generated data are too coarse for local
30 climate change impact assessments. To bridge the gap from the global to the local scale, Regional Climate Models have become a standard application (Jacob et al., 2014), using the output from GCMs as input or boundary conditions.

Although the information provided by both GCMs and RCMs is very valuable, both are biased with respect to the observations, especially for precipitation (Kotlarski et al., 2014). The biases can occur in any statistic and are commonly defined as “*a systematic difference between a simulated climate statistic and the corresponding real-world climate statistic*” (Maraun, 2016).
35 These biases are caused by temporal or spatial discretisation and unresolved or unrepresented physical processes (Teutschbein and Seibert, 2012; Cannon, 2016). An important example of the latter is convective precipitation, which can only be resolved by very high resolution models. Although the improvement of models is an important area of research (Prein et al., 2015; Kendon et al., 2017; Helsen et al., 2019; Fossier et al., 2020), such improved models are computationally expensive. As such, it is still necessary practice to statistically adapt the climate model output to adjust the biases (Christensen et al., 2008; Teutschbein and
40 Seibert, 2012; Maraun, 2016).

Many different bias-adjusting methods exist (Teutschbein and Seibert, 2012; Gutiérrez et al., 2019). They all calibrate a transfer function using the historical simulations and historical observations and apply this transfer function to the future simulations to generate future ‘observed values’ or an adjusted future. Of all the different methods, the quantile mapping method (Panofsky et al., 1958) was shown to be the generally best performing method (Rojas et al., 2011; Gudmundsson et al.,
45 2012). Quantile mapping adjusts biases in the full distribution, whereas most other methods only adjust biases in the mean and/or variance.

An important problem with quantile mapping and most other commonly used methods is that they are univariate and do not adjust biases in the multivariate correlation. Although quantile mapping can retain climate model multivariate correlation (Wilcke et al., 2013), the ability of univariate methods to improve the climate model’s multivariate correlation has been ques-
50 tioned (Hagemann et al., 2011; Ehret et al., 2012; Hewitson et al., 2014). This is important for impact assessment, as local impact models often need multiple input variables and many high-impact events are caused by the co-occurrence of multiple phenomena, the so-called ‘compound events’ (Zscheischler et al., 2018, 2020). For example, floods can be characterised by a rainfall-runoff model using evaporation and precipitation time series as an input. If the correlation between these variables is biased with respect to the observations, then it can be expected that the model output is biased as well. This results in a higher
55 uncertainty when using these models and thus in the resulting assessment. During the past decade, multiple methods have been developed to counter this problem. The first methods focused on the adjustment of two jointly occurring variables, most often precipitation and temperature, such as those by Piani and Haerter (2012) and Li et al. (2014). However, it became clear that ad-



justing only two variables would not suffice, hence many more methods have been developed that jointly adjust more variables, including those by Vrac and Friederichs (2015); Cannon (2016); Mehrotra and Sharma (2016); Dekens et al. (2017); Cannon (2018); Vrac (2018); Nguyen et al. (2018); Robin et al. (2019). Yet, the recent growth in availability of such methods comes along with a gap in the knowledge on their performance. In some studies, these methods have been compared with one or two older multivariate methods to reveal the improvements (Vrac and Friederichs, 2015; Cannon, 2018) or with univariate methods (Räty et al., 2018; Zscheischler et al., 2019; Meyer et al., 2019). Each of these three studies indicates that the univariate and multivariate methods lead to different results, yet it is difficult to conclude whether uni- or multivariate methods perform best. According to Zscheischler et al. (2019) multivariate methods have an added value. Räty et al. (2018) conclude that the multivariate methods and univariate methods performed similarly, while Meyer et al. (2019) could not draw definitive conclusions. These studies vary in set-up, adjusted variables and study area, which all could have caused the difference in added value. In all three studies, the same method, namely the Multivariate Bias Correction in n dimensions (MBCn) (Cannon, 2018) was the basis for comparison. Only recently, the first studies comparing multiple multivariate bias-adjusting methods were published (François et al., 2020; Guo et al., 2020). The study by François et al. (2020) focused on the different principles underlying the multivariate bias-adjusting methods and concluded that the choice of method should be based on the end user's goal. Besides, they also noticed that so far, all multivariate methods fail in representing the temporal structure of a time series. In contrast to the focus of François et al. (2020), Guo et al. (2020) studied the performance of multivariate bias-adjusting methods for climate change impact assessment and concluded that multivariate methods could be interesting in this context. However, they also noticed that the performance of the multivariate methods was lower in the more recent validation period and suggested that this could be caused by bias nonstationarity. As the use of multivariate bias-adjusting methods could be an important tool for climate change impact assessment, this deserves more attention.

The bias stationarity - or bias time invariance - assumption is the most important assumption for bias correction. It implies that the bias is the same in the calibration and validation or future periods and that the transfer function based on the calibration period can consequently be used in the future period. However, this assumption does not hold due to different types of nonstationarity induced by climate change, which may cause problems (Milly et al., 2008; Derbyshire, 2017). In the context of bias adjustment, this problem has been known for several years (Christensen et al., 2008; Ehret et al., 2012), but has not received a lot of attention. A few authors have tried to propose new types of bias relationships (Buser et al., 2009; Ho et al., 2012; Sunyer et al., 2014; Kerkhoff et al., 2014). Recently, it has been suggested that it is best to assume a non-monotonic bias change (Van Schaeybroeck and Vannitsem, 2016). Some authors suggested that bias nonstationarity could be an important source of uncertainty (Chen et al., 2015; Velázquez et al., 2015; Wang et al., 2018; Hui et al., 2019), but not all found clear indications of bias nonstationarity (Maraun, 2012; Piani et al., 2010; Maurer et al., 2013).

The availability of new methods and more data enables a more coherent assessment of the bias (non)stationarity issue. By comparing three bias-adjusting methods in a climate change context with possible bias nonstationarity, some of the remaining questions in François et al. (2020) and Guo et al. (2020) can be answered. The three multivariate bias-adjusting methods that will be compared in this study are 'Multivariate Recursive Quantile Nesting Bias Correction' (MRQNBC, Mehrotra and Sharma (2016)), MBCn (Cannon, 2018) and 'dynamical Optimal Transport Correction' (dOTC, Robin et al. (2019)). These three



methods give a broad view of the different multivariate bias adjustment principles, which we will elaborate on in Section 3.3. As a baseline, two univariate bias-adjusting methods will be used: Quantile Delta Mapping (QDM, Cannon et al. (2015)) and modified Quantile Delta Mapping (mQDM, Pham (2016)). QDM is a classical univariate bias-adjusting method and is chosen for this analysis as it is a robust and relatively common quantile mapping method, especially as one of the subroutines in the multivariate bias-adjusting methods (Mehrotra and Sharma, 2016; Nguyen et al., 2016; Cannon, 2018). mQDM, on the other hand, is one of the so-called ‘delta change’ methods, which are based on an adjustment of the historical time series. Using these univariate bias-adjusting methods, we can assess whether multivariate and univariate bias-adjusting methods differ in their response to possible bias nonstationarity.

The methods will be compared by applying them for the bias adjustment of precipitation, potential evaporation and temperature. The bias-adjusted time series will be used as inputs for a hydrological model in order to simulate the discharge. Discharge time series are the basis for flood hazard calculation, but can also be considered as an interesting source of validation themselves (Hakala et al., 2018). Although temperature is not needed as an input for the hydrological model, it is, together with precipitation, the most common variable to be adjusted in similar studies and therefore it is also included here. In order to mimic climate change context, the ‘historical’ or calibration time series runs from 1970 to 1989 and the ‘future’ or validation time series runs from 1998 to 2017, which is only recent past. In the latter time frame, effects of climate change are already visible (IPCC, 2013). The change of some biases from calibration to validation time series will be calculated, to indicate the extent of the bias nonstationarity. Maurer et al. (2013) proposed the R index for this purpose (see Section 2.4). Calculating the bias nonstationarity between both periods will give an indication of the impact of a changing bias on climate impact studies for the end of the 21st century. As Chen et al. (2015) mentioned: *“If biases are not constant over two very close time periods, there is little hope they will be stationary for periods separated by 50 to 100 years”*

2 Data and validation

2.1 Data

The observational data used were obtained from the Belgian Royal Meteorological Institute (RMI) Uccle observatory. The most important time series used is the 10-min precipitation amount, gauged with a Hellmann-Fuess pluviograph, from 1898 to 2018. An earlier version of this precipitation dataset was described by Demarée (2003) and analyzed in De Jongh et al. (2006). Multiple other studies have used this time series (Verhoest et al., 1997; Verstraeten et al., 2006; Vandenberghe et al., 2011; Willems, 2013). The 10-min precipitation time series was aggregated to daily level to be comparable with the other time series used.

For the multivariate method, the precipitation time series was combined with a 2 meter air temperature and potential evaporation time series. The daily potential evaporation was calculated by the RMI from 1901 to 2019, using the Penman formula for a grass reference surface (Penman, 1948) with variables measured at the Uccle observatory. Daily average temperatures were obtained using measurements from 1901 to 2019. As the last complete year for precipitation was 2017, the data were used from 1901 to 2017, amounting to 117 years of daily data.



The IPCC report (IPCC, 2013) clearly states the influence of climate change on different variables. For Belgium, this is illustrated by Fig. 1, in which the temperature and evaporation anomalies for the 21st century are all higher than the long-term mean value. However, for precipitation, the effect of climate change is not yet visible.

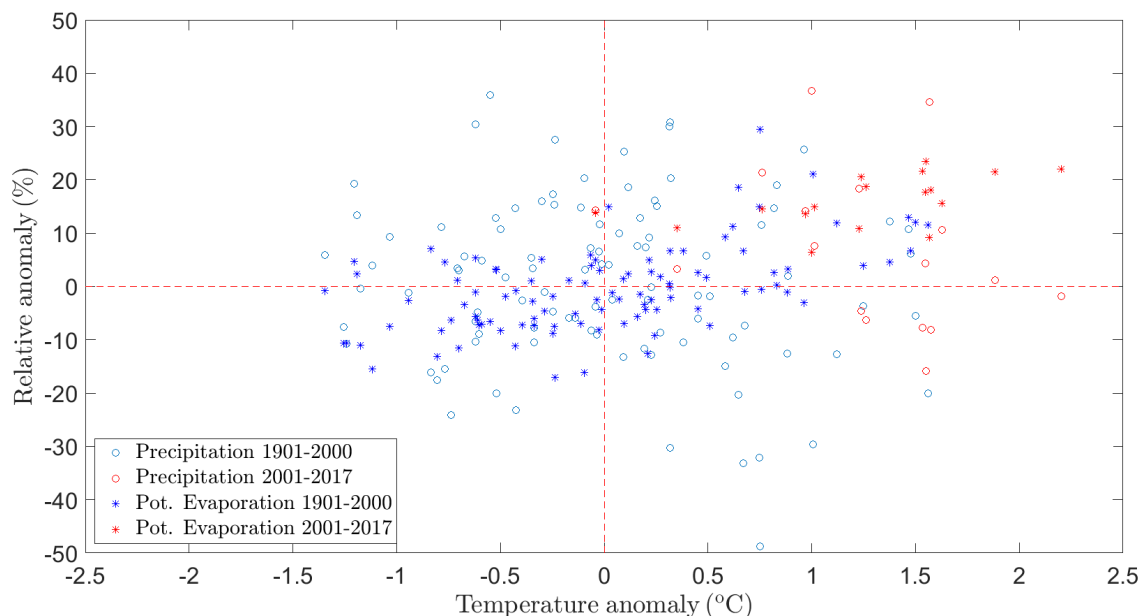


Figure 1. Yearly mean temperature, precipitation and evaporation anomalies for 1901-2017, compared with long-term mean value from 1920-1980. Red points are 21st century values.

For the simulations, data from the EURO-CORDEX project (Jacob et al., 2014) were used. The Rossby Centre regional climate model RCA4 was used (Strandberg et al., 2015) as it is one of the few RCMs with potential evaporation as an output variable. This RCM was forced with boundary conditions from the MPI-ESM-LR GCM (Popke et al., 2013). Historical data and scenario data for the grid cell comprising Uccle were respectively obtained for 1970-2005 and 2006-2100. The former time frame is limited by the earliest available data from the RCM. The latter time frame was only used until 2017, in accordance with the observational data. As climate change scenario, an RCP4.5 forcing was used in this paper (Van Vuuren et al., 2011). Since only ‘near future’ (from the model point of view) data were used, the choice of forcing does not have a large impact. However, when studying scenarios in a time frame further away from the present, using an ensemble of forcings is more relevant to be aware of the uncertainty regarding future climate change impact. Evaluations of the RCA4 model have shown that there is a bias in precipitation, especially in winter (Strandberg et al., 2015), but this bias is in line with the biases from other EURO-CORDEX models (Kotlarski et al., 2014).

As Uccle (near Brussels) is situated in a region with small topographic differences, it is assumed that the conditions in Uccle can be applied anywhere within the climate model grid cell and the variance within this cell is about the same. This assumption



can be made as long as the resulting adjusted data are not used for extremely localized studies, such as urban hydrology impact assessment.

2.2 Time frames

145 As mentioned in the introduction, it is important to assess bias-adjusting methods in a context they will be used in, i.e. under climate change conditions. The time series used in this study were chosen accordingly: 1970-1989 was chosen as the ‘historical’ or calibration time period and 1998-2017 was chosen as the ‘future’ or validation time period. Time series of 20 years were chosen here, although it is advised to use 30 years of data to have robust calculations (Berg et al., 2012; Reiter et al., 2018). However, as no climate model data prior to 1970 are available, using 30 years of data would have led to overlapping time series.

150 2.3 Validation framework

An important aspect in bias adjustment is the validation of the methods. Different methods are available, of which a pseudo-reality experiment (Maraun, 2012) is one of the most-used ones. In this method, each member of a model ensemble is in turn used as the reference in a cross-validation. However, while such a set-up is useful when comparing bias-adjustment methods, it only mimics a real application context. When sufficient observations are available, a ‘pseudo-projection’ setup
155 (Li et al., 2010) can be used. This set-up resembles a ‘differential split-sample testing’ (Klemeš, 1986) and is more in agreement with a practical application of bias-adjusting methods. Differential split-sample testing has been used in a bias adjustment context by Teutschbein and Seibert (2013), by constructing two time series with respectively the driest and wettest years. In our case study, it is assumed that the two time series differ enough because of climate change. Consequently, the approach is simple, and as the validation is not set in the future, it is considered a ‘pseudo-projection’.

160 Besides the choice of time frames and data, also the choice of validation indices is of key importance. Maraun and Widmann (2018a) stress that these indices should only be indirectly affected by the bias adjustment, as only validating on adjusted indices can be misleading. Such adjusted indices are the precipitation intensity, temperature and evaporation, which are used to build the transfer function in the historical setting and should be corrected by construction. Under bias stationarity, this correction will be carried over to the future, possibly hiding small inconsistencies that may arise for extreme values. If the bias is not
165 stationary, the effect might be different between adjusted and indirectly affected indices. As such, besides the three adjusted variables (indices 1 to 3 in Table 1) and their correlations (indices 4 to 12, which are directly adjusted by some of the methods), also indices based on the precipitation occurrence and on the discharge Q are used. The occurrence-based indices (13 to 16) allow for assessing how the methods influence the precipitation time series structure. The discharge-based indices (17 and 18) allow for the assessment of the impact of the different bias-adjusting methods on simulated river flow. The discharge-based
170 indices combine the information of the other indices by routing through the rainfall-runoff model. They are the most important aspect of the assessment, as they indicate the natural hazard. ETCCDI (Expert Team on Climate Change Detection and Indices) precipitation indices (Zhang et al., 2011) have also been considered and calculated. However, these are not included in this paper, as the differences in ETCCDI indices were minor and did not allow to clearly discern between the different methods. All



indices were calculated taking all days into account, instead of only calculating them on wet days, as some of the multivariate
 175 bias-adjusting methods do not discriminate between wet or dry days in their adjustment.

Table 1. Overview of the indices used

Nr	Index	Name
1	P_x	Precipitation amount percentile values, with x the percentile considered
2	T_x	Temperature percentile values, with x the percentile considered
3	E_x	Evaporation percentile values, with x the percentile considered
4	$\text{corr}_{P,E}$	Spearman correlation between the time series of P and E
5	$\text{corr}_{P,T}$	Spearman correlation between the time series of P and T
6	$\text{corr}_{E,T}$	Spearman correlation between the time series of E and T
7	$\text{crosscorr}_{P,E,0}$	Lag-0 crosscorrelation between the time series of P and E
8	$\text{crosscorr}_{P,T,0}$	Lag-0 crosscorrelation between the time series of P and T
9	$\text{crosscorr}_{E,T,0}$	Lag-0 crosscorrelation between the time series of E and T
10	$\text{crosscorr}_{P,E,1}$	Lag-1 crosscorrelation between the time series of P and E
11	$\text{crosscorr}_{P,T,1}$	Lag-1 crosscorrelation between the time series of P and T
12	$\text{crosscorr}_{E,T,1}$	Lag-1 crosscorrelation between the time series of E and T
13	P_{P00}	Precipitation transition probability from a dry to a dry day
14	P_{P10}	Precipitation transition probability from a wet to a dry day
15	N_{dry}	Number of dry days
16	$P_{\text{lag}1}$	Precipitation lag-1 auto-correlation
17	Q_x	Discharge percentiles, with x the percentile considered
18	Q_{T20}	20-year return period value of discharge

2.4 Bias nonstationarity

In a study on possible changes in bias, Maurer et al. (2013) proposed the R index:

$$R = 2 \frac{|\text{bias}_f - \text{bias}_h|}{|\text{bias}_f| + |\text{bias}_h|}, \quad (1)$$

where bias_f and bias_h are the biases in respectively the future and historical time series, calculated on the basis of the observa-
 180 tions and raw climate simulations. The R index takes a value between 0 and 2. If the index is greater than one, the difference
 in bias between the two sets is larger than the average bias of the model and it is likely that the bias adjustment would degrade
 the RCM output rather than improve it. The index is calculated for the indices used for validation in order to have an indication
 of the influence of bias nonstationarity on these indices. Besides for the indices, the R index is also calculated for the average
 and standard deviation of each variable, in order to be able to more easily visualise the changes in distribution.



185 2.5 Hydrological model

Similar to Pham et al. (2018), we use the Probability Distributed Model (PDM, Moore (2007)), a lumped conceptual rainfall-runoff model to calculate the discharge for the Grote Nete watershed in Belgium. This model uses precipitation and evaporation time series as inputs to generate a discharge time series. The PDM as used here was calibrated by Cabus (2008) using the Particle Swarm Optimization algorithm (PSO, Eberhart and Kennedy (1995)). As in Pham et al. (2018), it was assumed that the differences between meteorological conditions in the Grote Nete-watershed and Uccle are negligible, and that thus the adjusted data for the Uccle grid cell can be used as a forcing for the PDM. Furthermore, the goal is not to make predictions, but to assess the impact of different bias adjustment methods on the discharge values. To calculate the bias on the indices, observed, raw and adjusted RCM time series were used as forcing for this model. The discharge time series generated by the observations is considered to be the 'observed' discharge, and biases are calculated in comparison with this time series.

195 2.6 Validation metrics

The residual biases relative to the observations and to the model bias are often used in this paper to graphically present and interpret the results. These residual biases are based on the 'added value' concept (Di Luca et al., 2015) and enable a comparison based on two aspects. The first aspect is the performance in removing the bias, the second is the extent of the bias removal in comparison with the original value for the corresponding index for the observation time series. The use of the residual biases allows for a detailed study and comparison of the effect of bias adjustment on the different indices.

The residual bias relative to the observations RB_O for an index k is calculated as follows:

$$RB_O(k) = 1 - \frac{|\text{bias}_{\text{raw}(k)}| - |\text{bias}_{\text{adj}(k)}|}{|\text{obs}(k)|}, \quad (2)$$

with $\text{raw}(k)$ the raw climate model simulations, $\text{adj}(k)$ the adjusted climate model simulations and $\text{obs}(k)$ the observed values for index k .

The residual bias relative to the model bias RB_{MB} for an index k is calculated as follows:

$$RB_{MB}(k) = 1 - \frac{|\text{bias}_{\text{raw}(k)}| - |\text{bias}_{\text{adj}(k)}|}{|\text{bias}_{\text{raw}(k)}|}. \quad (3)$$

Absolute values are used in Eqs. (2) and (3) to compute the absolute difference between the raw and adjusted values, thus neglecting a possible change of sign of the bias. If the values of these residual biases are lower than 1 for an index, the method performs better than the raw RCM for this index. The best methods have low scores on both residual biases for as many indices as possible.

3 Bias-adjusting methods

3.1 Occurrence-bias adjustment: Thresholding

One of the deficiencies of RCMs, especially in Northwest Europe, are the so-called 'drizzle days' (Gutowski et al., 2003; Themeßl et al., 2012; Argüeso et al., 2013), i.e. the simulation of a small amount of precipitation on days that are supposed



215 to be dry. This has an influence on the temporal structure of the simulated time series and should thus be adjusted (Ines and
Hansen, 2006). This is commonly done in an occurrence-bias-adjusting step before the main step, the intensity-bias adjustment.
In this study, we use the thresholding occurrence-bias-adjusting method. Thresholding is one of the most common occurrence-
bias-adjusting methods and has been in use for many years (e.g. Hay and Clark (2003); Schmidli et al. (2006); Ines and Hansen
(2006)). This method is only applicable in regions where the assumption holds that the simulated time series has more wet
220 days than the observed time series. This is the case for Northwest Europe (Thiemeßl et al., 2012) and Belgium in particular. An
advanced version of the thresholding method is used here. To adjust the number of wet days, the frequencies of dry days in the
observations and in the simulations are calculated. The difference between the two frequencies, ΔN , is the number of days of
the simulated time series that have to be adapted. The simulated series is adapted by first sorting the wet days and thus only
changing the ΔN lowest days of the simulation time series by setting them to 0. ΔN is computed for the past and applied in
225 the future and consequently relies on the bias stationarity assumption. However, as thresholding is used prior to all methods,
the influence of possible bias nonstationarity on ΔN is assumed to be negligible.

In this advanced version of thresholding, some considerations are made. First, a day is considered wet if its simulated
precipitation amount is above 0.1 mm, to account for measurement errors in the observations. Second, the adjustment is done
on a monthly basis, to withhold a realistic temporal structure. This implies that to correct the number of wet days in month
230 m , all days of month m of the time series are selected. Third, both historical and future simulations are adjusted during the
same calculation step, to ensure a sound comparison during the intensity phase of the adjustment. If only either the historical
or future time series would have been adjusted, the assumption that the bias can be transferred from the historical to the future
time period would be impaired. The thresholding method is summarized in Algorithm 1.



Algorithm 1 Thresholding

Input:

Historical observations X^{ho}

Historical simulations X^{hs}

Future simulations X^{fs}

Output:

Adjusted historical $X_{\text{out}}^{\text{hs}}$ and future simulations $X_{\text{out}}^{\text{fs}}$

Initialization

for $m = 1 : 12$ **do**

Select data for month m : X_m^{ho} , X_m^{hs} and X_m^{fs} {Loop over months}

Calculate the percentage of dry days in month m for the historical observations

Calculate the percentage of dry days in month m for the historical simulations

Calculate ΔN for month m

{Adjustment of historical simulated time series}

Select and sort the wet days

Set the ΔN wet days with the lowest precipitation amount to 0

Restore the original order of the wet days of month m

Restore the full historical time series for month m

{Adjustment of future simulated time series}

Select and sort the wet days

Set the ΔN wet days with the lowest precipitation amount to 0

Restore the original order of the wet days of month m

Restore the full future time series for month m

{Reconstruction}

Replace the data in $X_{\text{out}}^{\text{hs}}$ with the adjusted data for month m

Replace the data in $X_{\text{out}}^{\text{fs}}$ with the adjusted data for month m

end for

3.2 Univariate intensity-bias-adjusting methods

235 3.2.1 Quantile Delta Mapping

The Quantile Delta Mapping (QDM) method was first proposed by Li et al. (2010). Its main idea is to preserve the climate simulation trends: it takes trend nonstationarity (changes in the simulated distribution) into account to a certain degree. Although it handles temperature adjustments well, it gives unrealistic values for precipitation and was therefore extended by Wang and



Chen (2014) for precipitation adjustment. A comparison with other quantile mapping methods by Cannon et al. (2015) showed
240 this method to perform best with respect to the preservation of trends. Cannon et al. (2015) bundled both the method by Li et al.
(2010) (*Equidistant CDF-matching*) and Wang and Chen (2014) (*Equiratio CDF-matching*) under the name *Quantile Delta
Mapping*, because of the similarity with delta change methods (which are described in e.g. Olsson et al. (2009), Willems and
Vrac (2011) and Rätty et al. (2014)).

Mathematically, this method can be written as

$$245 \quad x_i^{\text{fa}} = x_i^{\text{fs}} + F_{x^{\text{ho}}}^{-1}(F_{x^{\text{fs}}}(x^{\text{fs}})) - F_{x^{\text{hs}}}^{-1}(F_{x^{\text{fs}}}(x^{\text{fs}})) \quad (4)$$

in the additive case, and

$$x_i^{\text{fa}} = x_i^{\text{fs}} \frac{F_{x^{\text{ho}}}^{-1}(F_{x^{\text{fs}}}(x^{\text{fs}}))}{F_{x^{\text{hs}}}^{-1}(F_{x^{\text{fs}}}(x^{\text{fs}}))} \quad (5)$$

in the ratio or multiplicative case. The superscripts hs, ho, fs and fa indicate respectively the historical simulations, the
historical observations, the future simulations and the adjusted future. In this paper, the additive version is used for tem-
250 perature time series and the multiplicative one for precipitation and evaporation time series. This choice is based on the
work of Wang and Chen (2014), who have shown that using the additive adjustment for precipitation results in unrealistic
precipitation values and introduced a multiplicative adjustment. For evaporation, we follow the few available studies (e.g.
Lenderink et al. (2007)) in using the same adjustment as for precipitation.

To ensure the consistency of the time series, Themeßl et al. (2011) implemented a 61-day moving window. Here, a 91-day
255 moving window is opted for, as suggested by Rajczak et al. (2016) and Reiter et al. (2018). This enables the adjustment of
each day based on $91 \text{ days/year} \cdot 20 \text{ years} = 1820 \text{ days}$. These days were used to build an empirical CDF (as in Gudmundsson
et al. (2012); Gutjahr and Heinemann (2013), among others), because of the ease of application. The number of quantiles
implemented was determined automatically by the Matlab function `ecdf`, ranging between 200 and 400. It is also important
to note that for precipitation, Eq. (5) was applied only on the days considered wet, i.e. with a precipitation higher than 0.1 mm.
260 For consistency, a threshold of 0.1 mm was also used for evaporation. As no prior examples for evaporation adjustment were
available, we assumed this consistency was the best option. It is important to note that although QDM is only applied on wet
days, it can still transform low-precipitation wet days into days considered dry (e.g. with a precipitation amount $< 0.1 \text{ mm}$) if
the ratio in Eq. (5) is small enough.

3.2.2 Modified Quantile Delta Mapping

265 Pham (2016) proposed another version of QDM, following the delta change philosophy (Olsson et al., 2009; Willems and Vrac,
2011): the trend established by the RCM is assumed to be more thrust-worthy than the absolute value itself. When applying
this type of methods, the simulated change between the historical and the future is applied to the observations. Thus, instead
of the future simulations, the historical observations are adjusted to the future ‘observations’. As Johnson and Sharma (2011)
mention, this workflow could be problematic for future impact assessment, as it inherits the temporal structure of the historical
270 observations. This method is mathematically very similar to the QDM method, exchanging the roles of x^{fs} and x^{ho} . Thus, it is



named ‘modified Quantile Delta Mapping’ (mQDM), and can for the additive case be written as

$$x_i^{\text{fa}} = x_i^{\text{ho}} + F_{x_{\text{fs}}}^{-1} (F_{x_{\text{ho}}} (x^{\text{ho}})) - F_{x_{\text{hs}}}^{-1} (F_{x_{\text{ho}}} (x^{\text{ho}})). \quad (6)$$

The ratio version is mathematically written as

$$x_i^{\text{fa}} = x_i^{\text{ho}} \frac{F_{x_{\text{fs}}}^{-1} (F_{x_{\text{ho}}} (x^{\text{ho}}))}{F_{x_{\text{hs}}}^{-1} (F_{x_{\text{ho}}} (x^{\text{ho}}))}. \quad (7)$$

275 For the implementation, the same principles were used as for the QDM method: a 91-day moving window, empirical CDFs and a threshold of 0.1 mm/day to be considered as a wet day.

3.3 Multivariate intensity-bias-adjusting methods

The increasing number of multivariate bias-adjusting methods throughout the 2010s urges the need to classify them according to their properties. One possible classification was done by Vrac (2018), who proposed the ‘marginal/dependence’ versus the
280 ‘successive conditional’ approach. The former approach separately adjusts the 1D-marginal distributions and the dependence structure and is applied in e.g. Vrac and Friederichs (2015), Cannon (2018) and Vrac (2018). These two components are then recombined to obtain data that are close to the observations for both marginal and multivariate aspects. The latter approach consists of adjusting one given variable and then adjusting a second variable conditionally on the second variable: this procedure is applied successively to each variable. Examples can be found in e.g. Piani and Haerter (2012), Li et al. (2014) and
285 Dekens et al. (2017). According to Vrac (2018), the latter approach suffers from two main limitations. First, the adjustment is performed conditionally on the previously adjusted data. However, the adjustment is often applied in bins. As a result, for each variable, the amount of data available for each bin decreases, thus decreasing the robustness of the adjustment. Second, the ordering of the variables in the successive adjustments matters. For example, Li et al. (2014) point out that their ‘Joint Bias Correction for temperature’ (JBCt) and ‘Joint Bias Correction for precipitation’ (JBCp) methods, which respectively first ad-
290 just temperature and precipitation, differ in performance. For these two reasons, Vrac (2018) advocates for the use of the more robust and coherent ‘marginal/dependence’ approach. Robin et al. (2019) and François et al. (2020) extended this classification by introducing the all-in-one approach, which adjusts the marginal variables and the correlations simultaneously, ‘dynamical Optimal Transport Correction’ (dOTC) (Robin et al., 2019) being such a method.

Another perspective on the multivariate bias-adjusting methods is to consider the amount of temporal adjustment that is
295 allowed or applied by the method. This is important, as the amount of temporal adjustment is intrinsically linked with the main goal, the adjustment of the multivariate distribution of the variables. This distribution, in which the dependence is characterised by the underlying copula (Nelsen, 2006; Schölzel and Friederichs, 2008), can be estimated using the ranks. Thus, to adjust the multivariate distribution, the ranks of the climate model are replaced by those of the observations, using methods such as the ‘Schaake Shuffle’ (Clark et al., 2004; Vrac and Friederichs, 2015). This implies that the temporal structure and trends of
300 the climate model will be altered, which may have a considerable impact (François et al., 2020). This impact is especially large when multiday characteristics strongly matter, such as in applications as the hydrological example we use in this study (Addor and Seibert, 2014). Vrac (2018) mentions this necessity to modify the temporal structure and rank chronology of the



simulations. Yet, he also mentions that the extent of this modification is still a matter of debate. Cannon (2016) describes this as the ‘knobs’ that control whether marginal distributions, inter-variable or spatial dependence structure and temporal structure are more informed by the climate model or the observations. Thus, the choice between the temporal structure of the climate model and unbiased dependence structures is a trade-off that has to be made. Some methods, such as those by Vrac and Friederichs (2015), Mehrotra and Sharma (2016) and Nguyen et al. (2018) rely on the observations for their temporal properties, while other methods try to find the middle ground (e.g. Vrac (2018) and Cannon (2018)).

Our choice of multivariate bias-adjusting methods takes the above classification into account. The oldest method in the comparison is ‘Multivariate Recursive Quantile Nesting Bias Correction’ (Mehrotra and Sharma, 2016). This method completely replaces the simulated correlations by those of the observations and is a ‘marginal/dependence’ method according to François et al. (2020). ‘Multivariate Bias Correction in n dimensions’ (Cannon, 2018) is both a ‘marginal/dependence’ method and a method that tries to combine information from the climate model and the observations. The most recent method, ‘dynamical Optimal Transport Correction’ (Robin et al., 2019) differs considerably from the other two methods: it generalises the ‘transfer function’-principle using the ‘optimal transport’ paradigm (Villani, 2008), thereby defining a new category of multivariate bias-adjusting methods: the above-mentioned all-in-one approach. Although far from complete, by comparing these three methods, a broad view of the different approaches in multivariate bias adjustment can be obtained.

3.3.1 Multivariate Recursive Quantile Nesting Bias Correction

In 2016, Mehrotra and Sharma proposed a new multivariate bias adjustment method, named ‘Multivariate Recursive Quantile Nesting Bias Correction’ (MRQNBC), based on a combination of several older methods by Johnson and Sharma (2012), Mehrotra and Sharma (2012) and Mehrotra and Sharma (2015). The underlying idea of this method is to adjust on more than one timescale, an idea that most bias-adjusting methods do not incorporate (Haerter et al., 2011). This adjustment on multiple timescales is applied by adjusting the biases in lag-0- and lag-1-auto and the cross-correlation coefficients, i.e. the persistence attributes, instead of focusing on the mean or the distribution.

As a first step in this method, QDM is applied separately on variables to adjust the empirical CDFs. This is followed by a multivariate bias adjustment to adjust the lag-0 and lag-1 auto and cross-correlation coefficients. This combination of univariate and multivariate bias-adjusting methods is applied on all time scales. For the multivariate adjustment, two models are used: a multivariate first-order autoregressive (AR(1)) model with constant parameters at the daily and yearly level, and a multivariate AR(1) model with periodic parameters (Salas, 1980) at the monthly and seasonal level. All steps are applied to the different types of data: historical observations of temperature, evaporation and precipitation (combined in the matrix \mathbf{X}^{ho}), historical climate model simulations of the three variables (the matrix \mathbf{X}^{hs}) and climate model projections of the three variables (the matrix \mathbf{X}^{fs}), which have to be adjusted. All these datasets are of size $T \times N$, with T the number of time steps and N the number of variables.

The quantile-mapped future GCM time series for time step t is denoted as \mathbf{X}_t^{fa} . The standardised versions of this time series and of the observed time series are denoted as $\check{\mathbf{X}}_t^{\text{fa}}$ and $\check{\mathbf{X}}_t^{\text{ho}}$, respectively. Using the standardised time (zero mean and unit variance) series, the multivariate AR(1) model with constant parameters (MAR) for the observed and GCM multivariate time



series can be expressed as (Mehrotra and Sharma, 2016):

$$\check{\mathbf{X}}_t^{\text{ho}} = \mathbf{C}\check{\mathbf{X}}_{t-1}^{\text{ho}} + \mathbf{D}\epsilon_t, \quad (8)$$

and

$$340 \quad \check{\mathbf{X}}_t^{\text{fa}} = \mathbf{E}\check{\mathbf{X}}_{t-1}^{\text{fa}} + \mathbf{F}\epsilon_t, \quad (9)$$

with \mathbf{C} and \mathbf{D} the coefficient matrices of $\check{\mathbf{X}}_t^{\text{ho}}$, \mathbf{E} and \mathbf{F} the coefficient matrices of $\check{\mathbf{X}}_t^{\text{fa}}$ and ϵ_t a white noise term. The coefficient matrices are calculated using the $N \times N$ lag-0 and lag-1 cross-correlation matrices \mathbf{M}_0 and \mathbf{M}_1 . Using the standardised time series, the elements of these matrices can be expressed as (Salas, 1980):

$$m_0^{i,j} = \frac{1}{T} \sum_{t=1}^T x_t^i x_t^j \quad (10a)$$

$$345 \quad m_1^{i,j} = \frac{1}{T-1} \sum_{t=1}^T x_t^i x_{t+1}^j, \quad (10b)$$

with i and j the column numbers of $\check{\mathbf{X}}_t$, referring to the variables whose correlation is calculated. This enables the calculation of \mathbf{C} and \mathbf{E} as (Matalas, 1967):

$$\mathbf{C} = \mathbf{M}_1^{\text{ho}} \mathbf{M}_0^{\text{ho}-1}, \quad (11a)$$

$$\mathbf{E} = \mathbf{M}_1^{\text{fa}} \mathbf{M}_0^{\text{fa}-1}, \quad (11b)$$

350 and of \mathbf{D} and \mathbf{F} via

$$\mathbf{D}\mathbf{D}^T = \mathbf{M}_0^{\text{ho}} - \mathbf{M}_1^{\text{ho}} \mathbf{M}_0^{\text{ho}-1} \mathbf{M}_1^{\text{ho}T} \quad (12a)$$

$$\mathbf{F}\mathbf{F}^T = \mathbf{M}_0^{\text{fa}} - \mathbf{M}_1^{\text{fa}} \mathbf{M}_0^{\text{fa}-1} \mathbf{M}_1^{\text{fa}T}, \quad (12b)$$

which can be solved using the eigenvalues and eigenvectors of $\mathbf{D}\mathbf{D}^T$ or $\mathbf{F}\mathbf{F}^T$:

$$\mathbf{D} = \mathbf{V}\sqrt{\mathbf{S}}\mathbf{V}^T, \quad (13a)$$

$$355 \quad \mathbf{F} = \mathbf{V}\sqrt{\mathbf{S}}\mathbf{V}^T, \quad (13b)$$

with \mathbf{V} the matrix of eigenvectors and \mathbf{S} a diagonal matrix with the corresponding eigenvalues.

The multivariate bias adjustment is then implemented by removing the lag-0 and lag-1 auto- and cross-correlations from the future time series $\check{\mathbf{X}}_t^{\text{fa}}$ (the matrices \mathbf{E} and \mathbf{F}) and applying the observed lag-0 and lag-1 auto- and cross-correlations (\mathbf{C} and \mathbf{D}) to the future time series and thus creating a modified future time series, $\check{\mathbf{X}}_t^{\prime\text{fa}}$. These steps are applied by first rearranging

360 and simplifying Eq. (9) for ϵ_t :

$$\epsilon_t = \mathbf{F}^{-1} \left(\check{\mathbf{X}}_t^{\text{fa}} - \mathbf{E}\check{\mathbf{X}}_{t-1}^{\text{fa}} \right), \quad (14)$$



with ϵ_t now a standardised vector of N variables calculated by removing the lag-0 and lag-1 auto- and cross-correlations from the $\check{\check{X}}_t^{\text{fa}}$ time series. This vector is plugged into Eq. (8) along with the matrices \mathbf{C} and \mathbf{D} in which $\check{\check{X}}_t^{\text{fa}}$ is used instead of $\check{\check{X}}_t^{\text{ho}}$ to obtain the modified time series:

$$365 \quad \check{\check{X}}_t^{\text{fa}} = \mathbf{C}\check{\check{X}}_{t-1}^{\text{fa}} + \mathbf{D}\mathbf{F}^{-1} \left(\check{\check{X}}_t^{\text{fa}} - \mathbf{E}\check{\check{X}}_{t-1}^{\text{fa}} \right), \quad (15)$$

which can be rearranged as:

$$\check{\check{X}}_t^{\text{fa}} = \mathbf{C}\check{\check{X}}_{t-1}^{\text{fa}} + \mathbf{D}\mathbf{F}^{-1} \check{\check{X}}_t^{\text{fa}} - \mathbf{D}\mathbf{F}^{-1} \mathbf{E}\check{\check{X}}_{t-1}^{\text{fa}}. \quad (16)$$

This model preserves the observed persistence attributes. As a last step, destandardising results in the bias-adjusted time series \mathbf{X}_t^{fa} .

370 When using the multivariate AR(1) with periodic parameters (PMAR), the parameters are derived separately for each period to allow for periodicity. In this case, the vectors $\mathbf{X}_{t,\tau}^{\text{ho}}$ and $\mathbf{X}_{t,\tau}^{\text{fa}}$ respectively represent the observed and quantile mapped GCM time series. The subscript t refers to the year and the subscript τ to a specific period in the year.

The elements of the periodic version of \mathbf{M}_0 and \mathbf{M}_1 can be calculated as (Salas, 1980):

$$m_{0,\tau}^{i,j} = \frac{\sum_{t=1}^{T_\tau} (x_{t,\tau}^i - \bar{x}_\tau^i) (x_{t,\tau}^j - \bar{x}_\tau^j)}{T_\tau s_\tau^i s_\tau^j} \quad (17a)$$

$$375 \quad m_{1,\tau}^{i,j} = \frac{\sum_{t=1}^{T_\tau} (x_{t,\tau}^i - \bar{x}_\tau^i) (x_{t,\tau-1}^j - \bar{x}_{\tau-1}^j)}{T_\tau s_\tau^i s_{\tau-1}^j}, \quad (17b)$$

with T_τ the number of time steps of the period τ , \bar{x}_τ and $\bar{x}_{\tau-1}$ the mean of periods τ and $\tau - 1$ (for instance, if τ is summer, than $\tau - 1$ is spring) and s_τ and $s_{\tau-1}$ the standard deviations of periods τ and $\tau - 1$. The correlation matrices are calculated in the same way as in the non-periodic steps. The only difference is that they are calculated for every period (e.g. separately for every season or month). For every time step in period τ , the corresponding value can be adjusted as follows to preserve the
 380 observed persistence attributes:

$$\check{\check{X}}_{t,\tau}^{\text{fa}} = \mathbf{C}_\tau \check{\check{X}}_{t,\tau-1}^{\text{fa}} + \mathbf{D}_\tau \mathbf{F}_\tau^{-1} \check{\check{X}}_{t,\tau}^{\text{fa}} - \mathbf{D}_\tau \mathbf{F}_\tau^{-1} \mathbf{E}_\tau \check{\check{X}}_{t,\tau-1}^{\text{fa}}. \quad (18)$$

The different time steps are combined with the nesting method proposed in Johnson and Sharma (2012) and Mehrotra and Sharma (2015). First, QDM (as described in Section 3.2.1) is applied at a daily level, followed by MAR. These adjusted time series are then aggregated and averaged to form a monthly time series, which is adjusted by QDM, standardised and adjusted
 385 by PMAR. Note that the standardisation of the aggregated time series does not imply that the variables of a period τ of that time series have zero mean and unit variance. The results of the monthly adjustment are aggregated and averaged to form seasonal time series, which are also adjusted using QDM, standardised and adjusted by PMAR. As a last nesting step, the results are once more aggregated and averaged to build an annual time series, which is adjusted using QDM and MAR. The outcomes of



all these steps are combined into a weighting factor that is used to modify the daily time series accordingly (Srikanthan and Pegram, 2009):

$$\mathbf{X}_{t,j,s,i}^{//fa} = \left(\frac{\mathbf{Y}_{j,s,i}^{/fa}}{\mathbf{Y}_{j,s,i}^{fa}} \right) \left(\frac{\mathbf{Z}_{s,i}^{/fa}}{\mathbf{Z}_{s,i}^{fa}} \right) \left(\frac{\mathbf{A}_i^{/fa}}{\mathbf{A}_i^{fa}} \right) \mathbf{X}_{t,j,s,i}^{fa}, \quad (19)$$

with t the day, j the month, s the season, i the year, $\mathbf{Y}_{j,s,i}^{/fa}$ the monthly adjusted value, $\mathbf{Y}_{j,s,i}^{fa}$ the aggregated-averaged monthly value, $\mathbf{Z}_{s,i}^{/fa}$ the seasonal adjusted value, $\mathbf{Z}_{s,i}^{fa}$ the aggregated-averaged seasonal value, $\mathbf{A}_i^{/fa}$ the adjusted yearly value and \mathbf{A}_i^{fa} the aggregated-averaged yearly value. The full procedure is summarised in Algorithm 2.

Algorithm 2 MRQNBC

Input:

- Daily historical observations \mathbf{X}^{ho}
- Daily historical simulations \mathbf{X}^{hs}
- Daily future simulations \mathbf{X}^{fs}

Output:

- Adjusted future simulations $\mathbf{X}^{//fa}$

for #Timescales do

- Apply QDM to calculate \mathbf{X}^{fa}
- Standardise \mathbf{X}^{ho} and \mathbf{X}^{fa}
- Calculate matrices **C** and **D** of $\check{\mathbf{X}}^{\text{ho}}$ and **E** and **F** of $\check{\mathbf{X}}^{\text{fa}}$ {The calculations of **C**, **D**, **E** and **F** depend on the periodicity of the timescale}
- Apply the persistence adjustment to calculate $\check{\check{\mathbf{X}}}^{/fa}$
- Destandardise $\check{\check{\mathbf{X}}}^{/fa}$
- Aggregate \mathbf{X}^{ho} and \mathbf{X}^{fa} to the higher timescale {Except for the yearly timescale}

end for

- Calculate weighting factors for all timescales except the daily timescale
 - Calculate the final adjusted daily value $\mathbf{X}^{//fa}$
-

395 The nesting method cannot fully remove biases at all time scales, thus Mehrotra and Sharma (2016) suggested to repeat the complete procedure multiple times. However, in our case this seemed to exacerbate the results, so the method was run only once.



3.3.2 Multivariate Bias Correction in n dimensions

In 2018, Cannon (2018) proposed the ‘Multivariate Bias correction in n dimensions’ (MBCn) method as a flexible multivariate
 400 bias-adjusting method. The method’s flexibility has attracted some attention, as it has already been used in multiple studies
 (Räty et al., 2018; Zscheischler et al., 2019; Meyer et al., 2019; François et al., 2020). This method consists of three steps.
 First, the multivariate data are rotated using a randomly generated orthogonal rotation matrix, adjusted with the additive form
 of QDM, and rotated back until the calibration period model simulations converge to the observations. This convergence is
 verified on the basis of the energy distance (Rizzo and Székely, 2016). Second, the validation period simulations are adjusted
 405 using QDM, as this method preserves the simulated trends. As the last step, these adjusted time series are shuffled using the
 Schaake Shuffle (Clark et al., 2004) based on the rank order of the rotated dataset.

Considering the j -th iteration of the method, denoted by the subscript $[j]$, the first step consists of rotating the data sets
 using an $N \times N$ randomly generated orthogonal rotation matrix $\mathbf{R}_{[j]}$. This orthogonal rotation matrix was created using the
 algorithm by Mezzadri (2007, pg. 597). This rotation is formulated as

$$\begin{aligned}
 \tilde{\mathbf{X}}_{[j]}^{\text{hs}} &= \mathbf{X}_{[j]}^{\text{hs}} \mathbf{R}_{[j]} \\
 410 \quad \tilde{\mathbf{X}}_{[j]}^{\text{fs}} &= \mathbf{X}_{[j]}^{\text{fs}} \mathbf{R}_{[j]} \\
 \tilde{\mathbf{X}}_{[j]}^{\text{ho}} &= \mathbf{X}_{[j]}^{\text{ho}} \mathbf{R}_{[j]}
 \end{aligned} \tag{20}$$

with $\tilde{\mathbf{X}}_{[j]}$ the resulting rotated matrix. In the next step, additive quantile delta mapping is applied to each variable in $\tilde{\mathbf{X}}_{[j]}^{\text{hs}}$ and
 $\tilde{\mathbf{X}}_{[j]}^{\text{fs}}$, using the corresponding variable in $\tilde{\mathbf{X}}_{[j]}^{\text{ho}}$ as the target. The resulting matrices $\mathbf{X}_{[j]}^{\text{hc}}$ and $\mathbf{X}_{[j]}^{\text{fa}}$ are rotated back:

$$\begin{aligned}
 \mathbf{X}_{[j+1]}^{\text{hs}} &= \mathbf{X}_{[j]}^{\text{hc}} \mathbf{R}_{[j]}^{-1} \\
 \mathbf{X}_{[j+1]}^{\text{fs}} &= \mathbf{X}_{[j]}^{\text{fa}} \mathbf{R}_{[j]}^{-1} \\
 \mathbf{X}_{[j+1]}^{\text{ho}} &= \mathbf{X}_{[j]}^{\text{ho}}
 \end{aligned} \tag{21}$$

These steps have to be repeated until the multivariate distribution of $\mathbf{X}_{[j+1]}^{\text{hs}}$ matches \mathbf{X}^{ho} . The similarity is measured using
 415 the (squared) energy distance (Székely and Rizzo, 2004, 2013; Rizzo and Székely, 2016), a measure of statistical discrepancy
 between two multivariate distributions. For two N -dimensional independent random vectors \mathbf{x} and \mathbf{y} with respective CDFs F
 and G , this measure is given by:

$$D^2(F, G) = 2E[\|\mathbf{x} - \mathbf{y}\|] - E[\|\mathbf{x} - \mathbf{x}'\|] - E[\|\mathbf{y} - \mathbf{y}'\|] \geq 0 \tag{22}$$



with E the expected value, $\| \cdot \|$ the Euclidean norm and \mathbf{x}' and \mathbf{y}' and \mathbf{x} and \mathbf{y} i.i.d.. A practical way of calculating this
 420 measure is as follows (Székely and Rizzo, 2013), with $\mathbf{X} = (X_1, \dots, X_{n_1})$ and $\mathbf{Y} = (Y_1, \dots, Y_{n_2})$:

$$\begin{aligned}
 D(\mathbf{X}, \mathbf{Y}) &= \frac{2}{n_1 n_2} \sum_{i=1}^{n_1} \sum_{m=1}^{n_2} |X_i - Y_m| \\
 &\quad - \frac{1}{n_1^2} \sum_{i=1}^{n_1} \sum_{j=1}^{n_1} |X_i - X_j| \\
 &\quad - \frac{1}{n_2^2} \sum_{l=1}^{n_2} \sum_{m=1}^{n_2} |Y_l - Y_m|,
 \end{aligned} \tag{23}$$

with i, j, l and m denoting the time steps.

As a last step, the preservation of trends of QDM has to be combined with the restoration of the multivariate ranks by
 the transformations. To do this, first either the additive or multiplicative version of quantile delta mapping (depending on
 425 the variable) has to be applied to each variable of the original data set \mathbf{X}^{fs} , using \mathbf{X}^{hs} and \mathbf{X}^{ho} as historical baseline data.
 As a final step, the elements of each column of \mathbf{X}^{fa} are reordered following a method known as the ‘Schaake Shuffle’
 (Clark et al., 2004; Vrac and Friederichs, 2015). In this method, the ranks of each variable of \mathbf{X}^{fa} are swapped with the ranks
 of the corresponding variables of $\mathbf{X}_{[j+1]}^{\text{fs}}$, thus reordering the time series’ structure according to the ranks of the observations.
 The Schaake Shuffle can be mathematically formulated as follows (Clark et al., 2004). Let \mathbf{X} be a vector of n time steps of a
 430 variable, and χ be the sorted vector of \mathbf{X} , that is:

$$\mathbf{X} = (x_1, x_2, \dots, x_n), \text{ and} \tag{24}$$

$$\chi = (x_{(1)}, x_{(2)}, \dots, x_{(n)}), \quad x_{(1)} \leq x_{(2)} \dots \leq x_{(n)} \tag{25}$$

Let \mathbf{Y} be a vector of n historical observations and γ be the sorted vector of \mathbf{Y} :

$$\mathbf{Y} = (y_1, y_2, \dots, y_n), \text{ and} \tag{26}$$

$$435 \quad \gamma = (y_{(1)}, y_{(2)}, \dots, y_{(n)}), \quad y_{(1)} \leq y_{(2)} \dots \leq y_{(n)}. \tag{27}$$

\mathbf{B} is then the vector of indices describing the original observation number that corresponds to the values in the ordered vector γ .
 The main step of the Schaake Shuffle is to reconstruct the reordered vector \mathbf{X}^{ss} :

$$\mathbf{X}^{\text{ss}} = (x_1^{\text{ss}}, x_2^{\text{ss}}, \dots, x_n^{\text{ss}}) \tag{28}$$

where

$$440 \quad x_q^{\text{ss}} = x_{(r)}, \quad q = \mathbf{B}[r], \text{ and } r = 1, \dots, n. \tag{29}$$

In MBCn, \mathbf{X}^{fa} is sorted according to $\mathbf{X}_{[j+1]}^{\text{fs}}$ instead of respectively \mathbf{X} and \mathbf{Y} , resulting in a final adjusted data set $\mathbf{X}^{\text{fa}'}$. To
 account for ties in this procedure, a small random value is added before calculating the ranks of $\mathbf{X}_{[j+1]}^{\text{fs}}$. The version of QDM



applied in the step prior to the Schaake Shuffle is the same as in Section 3.2.1. As such, the shuffling procedure is the only difference with the univariate bias-adjustment, implying that differences in performance can be related to it.

445 The MBCn method was shown by Cannon (2018) to outperform many earlier multivariate bias-adjusting methods, such as
the EC-BC (Vrac and Friederichs, 2015), the JBC (Li et al., 2014), MBCr and MBCp methods (Cannon, 2016). In contrast,
Cannon (2018) also pointed out some problems. Depending on the number of variables, the computational cost can get too
high and convergence speed too low, or overfitting might become an issue. The first problem can be tackled by implementing
sufficient time steps when having a lot of variables. Second, to address the convergence speed, Pitié et al. (2007) suggested
450 using a deterministic selection of rotation matrices that maximizes the distance between rotation axis sets instead of randomly
generating them. It is also suggested by Cannon (2018) to use the most efficient form of quantile mapping and to limit the
use of an advanced quantile mapping method to the last step. Third, to avoid overfitting (and to reduce the computational
cost), early stopping is also suggested by Cannon (2018) (e.g. Prechelt (1998)). As the number of variables is limited in this
case, overfitting did not seem to be a problem. Yet, to reduce unnecessary computational costs, the similarity in consecutive
455 energy distances was used as a measure to stop the computation. A tolerance of 0.0001 was used: if the difference between
two consecutively calculated energy distances was lower, the computation was halted. The full procedure is summarised in
Algorithm 3.

Algorithm 3 MBCn

Input:

Historical observations \mathbf{X}^{ho}

Historical simulations \mathbf{X}^{hs}

Future simulations \mathbf{X}^{fs}

Output:

Adjusted future simulations \mathbf{X}^{fa}

Initialisation: tolerance ϵ and initial energy distance difference ΔD_0

while $\Delta D > \epsilon$ **do**

Randomly generate a rotation matrix \mathbf{R}

Rotate \mathbf{X}^{ho} , \mathbf{X}^{hs} and \mathbf{X}^{fs}

Apply the additive form of QDM

Rotate \mathbf{X}^{ho} , \mathbf{X}^{hs} and \mathbf{X}^{fs} back

Calculate the energy distance D between \mathbf{X}^{hs} and \mathbf{X}^{ho}

Calculate the decrease in energy distance ΔD

end while

Apply QDM to the original inputs to calculate \mathbf{X}^{fa}

Apply the Schaake Shuffle based on the rotated future simulations to calculate \mathbf{X}^{fa}



3.3.3 Dynamical Optimal Transport Correction

460 Recently, Robin et al. (2019) indicated that the notion of a transfer function in quantile mapping can be generalised to the theory of optimal transport. Optimal transport is a way to measure the dissimilarity between two probability distributions and to use this as a means for transforming the distributions in the most optimal way (Villani, 2008; Peyré and Cuturi, 2019).

Optimal transport was used by Robin et al. (2019) to adjust the bias of a multivariate data set in the ‘dynamical Optimal Transport Correction’ method (dOTC), which extends the ‘CDF-transform’ (CDF-t) bias-adjusting method (Michelangeli et al., 2009). dOTC calculates the optimal transport plans from \mathbf{X}^{ho} to \mathbf{X}^{hs} (the bias between the model and the simulations) and 465 from \mathbf{X}^{hs} to \mathbf{X}^{fs} (the evolution of the model). The combination of both optimal transport plans allows for bias adjustment while preserving the trend of the model.

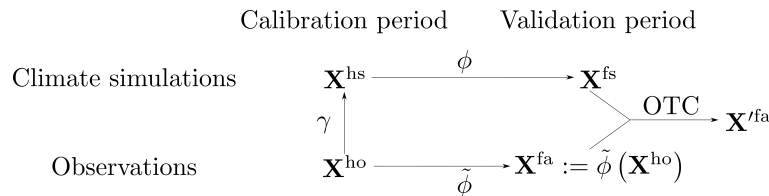


Figure 2. The different transport plans and transformations used in dOTC. Based on Robin et al. (2019).

Optimal transport is applied on the basis of an optimal transport plan. The optimal plan between \mathbf{X}^{ho} and \mathbf{X}^{hs} is denoted as γ . The second optimal plan between \mathbf{X}^{hs} and \mathbf{X}^{fs} is denoted as ϕ . The goal is to transform ϕ according to γ , defining a new plan $\tilde{\phi}$. This optimal plan estimates $\mathbf{X}^{\text{fa}} = \tilde{\phi}(\mathbf{X}^{\text{ho}})$. Finally, \mathbf{X}^{fs} is adjusted with respect to \mathbf{X}^{fa} , creating the final adjusted 470 $\mathbf{X}^{\text{fa-tilde}}$. These steps are summarised in Fig. 2.

For the definition of the optimal plan, the ‘Optimal Transport Correction’ (OTC) (Robin et al., 2019) is used. First, the empirical distributions $\hat{\mathbf{P}}_{\mathbf{X}^{\text{ho}}}$ and $\hat{\mathbf{P}}_{\mathbf{X}^{\text{hs}}}$ have to be calculated. To achieve this, the subspace of \mathbb{R}^N that contains the data is partitioned into regularly spaced cells, generically denoted \mathbf{c}_i^* , with N the number of variables of \mathbf{X}^{ho} and \mathbf{X}^{hs} . Using this notation, $\hat{\mathbf{P}}_{\mathbf{X}^{\text{ho}}}$ and $\hat{\mathbf{P}}_{\mathbf{X}^{\text{hs}}}$ can be estimated using the relative frequencies $p_{\mathbf{X}^{\text{ho}}}$ and $p_{\mathbf{X}^{\text{hs}}}$ as:

$$475 \quad p_{\mathbf{X}^{\text{ho}},i} = \frac{1}{m} \sum_{k=1}^m \mathbf{1}(\mathbf{X}_k^{\text{ho}} \in \mathbf{c}_i^*), \quad (30a)$$

$$p_{\mathbf{X}^{\text{hs}},i} = \frac{1}{n} \sum_{l=1}^n \mathbf{1}(\mathbf{X}_l^{\text{hs}} \in \mathbf{c}_i^*), \quad (30b)$$

with $\mathbf{1}$ the indicator function and m and n the total number of time steps of respectively \mathbf{X}^{ho} and \mathbf{X}^{hs} . Thus, the distributions are essentially estimated by counting the number of observations of each time series within each cell. The optimal plan γ between \mathbf{X}^{ho} and \mathbf{X}^{hs} can be estimated as:

$$480 \quad \hat{\gamma} = \sum_{i,j=1}^{I,J} \gamma_{i,j}. \quad (31)$$



The coefficients $\gamma_{i,j}$ are the probabilities to transform an observation of \mathbf{X}^{ho} in cell \mathbf{c}_i^* into an observation of \mathbf{X}^{hs} in cell \mathbf{c}_j^* and I and J are the total number of cells containing observations of respectively \mathbf{X}^{ho} and \mathbf{X}^{hs} . Note that from here on, \mathbf{c}_i^* and \mathbf{c}_j^* denote only those cells containing observations of respectively \mathbf{X}^{ho} and \mathbf{X}^{hs} and that ‘observation’ is used generically for both observed and simulated time series. The coefficients are unknown, but obey the marginal properties:

$$485 \quad \sum_{j=1}^J \gamma_{i,j} = p_{\mathbf{X}^{\text{ho}},i} \quad (32a)$$

and

$$\sum_{i=1}^I \gamma_{i,j} = p_{\mathbf{X}^{\text{hs}},j}. \quad (32b)$$

Central to the optimal transport theorem is the cost function C (Villani, 2008), which can here be approximated by

$$\hat{C}(\hat{\gamma}) = \sum_{i,j=1}^{I,J} \|\mathbf{c}_i - \mathbf{c}_j\|^2 \gamma_{i,j}, \quad (33)$$

with $\|\cdot\|$ the Euclidean norm and \mathbf{c}_i and \mathbf{c}_j the centres of the cells defined above. Finding $\gamma_{i,j}$ comes down to solving the
 490 problem defined by the constraints of Eqs. (32) and minimisation of Eq. (33). Here, Sinkhorn’s algorithm (Cuturi, 2013) is used to find the solution to this problem and thus the optimal plan γ . Using this optimal plan, a vector of probabilities with length J can be defined for each cell \mathbf{c}_i^* of \mathbf{X}^{ho} , each element j corresponding to the probability that an observation in that cell \mathbf{c}_i^* will be transformed into an observation in cell \mathbf{c}_j^* of \mathbf{X}^{hs} . In the transformation, the vectors of probabilities are used to introduce stochasticity, by sampling from these vectors the element j corresponding to cell \mathbf{c}_j^* . The stochastic transformation
 495 of an observation of \mathbf{X}^{ho} into an observation of \mathbf{X}^{hs} can be repeated to create an ensemble of results. This ensemble accounts for random weather effects and can thus be considered to be more similar to the true range of observations.

The optimal plan ϕ can be calculated analogously. This optimal plan ϕ transforms an observation of \mathbf{X}^{hs} in cell \mathbf{c}_j^* into an observation of \mathbf{X}^{fs} in cell \mathbf{c}_k^* , with \mathbf{c}_k^* defined analogously to \mathbf{c}_i^* and \mathbf{c}_j^* . What distinguishes dOTC from OTC is the next phase, in which ϕ is transformed according to γ , resulting in $\tilde{\phi}$. This is conducted in three steps, the first being is the transformation of
 500 ϕ into a vector. The vector $\mathbf{v}_{jk} := \mathbf{c}_k - \mathbf{c}_j$ represents the climatic trend from an observation of \mathbf{X}^{hs} in cell \mathbf{c}_j^* to an observation of \mathbf{X}^{fs} in cell \mathbf{c}_k^* . The second step is the transfer according to γ . The result $\tilde{\phi}$ can be defined by translating the observations of \mathbf{X}^{ho} along their respective vectors \mathbf{v}_{jk} : an observation of \mathbf{X}^{fa} is then given by $\mathbf{X}_t^{\text{ho}} + \mathbf{v}_{jk}$, with \mathbf{X}_t^{ho} the observation of \mathbf{X}^{ho} at time step t . However, the translation of \mathbf{X}_t^{ho} along vector \mathbf{v}_{jk} does not always define an optimal transport plan: the vector has to be adapted to \mathbf{X}^{ho} , which is the third step. In this step, a matrix factor \mathbf{D} is introduced, which rescales the vector \mathbf{v}_{jk} . This
 505 rescaling is actually the replacement of the scale of \mathbf{X}^{hs} by that of \mathbf{X}^{ho} . A Cholesky decomposition of the covariance matrix has been proposed for this rescaling (Bárdossy and Pegram, 2012; Cannon, 2016). Denoting the covariance matrix as Σ , and the Cholesky decomposition as $\text{Cho}(\Sigma)$, Robin et al. (2019) proposed to multiply \mathbf{v}_{jk} by the following matrix:

$$\mathbf{D} := \text{Cho}(\Sigma_{\mathbf{X}^{\text{ho}}}) \cdot \text{Cho}(\Sigma_{\mathbf{X}^{\text{hs}}})^{-1}. \quad (34)$$



Robin et al. (2019) remark that the Cholesky decomposition only exists if Σ is symmetric and positive-definite. Some covari-
 510 ance matrices, such as those of highly correlated random variables, do not have this property. Σ must then be slightly perturbed
 to be positive-definite (Higham, 1988; Knol and ten Berge, 1989). It is also possible for the Cholesky decomposition to be
 poorly estimated if the available data are too small compared to the dimension. In that case, it is suggested to replace the matrix
D by the diagonal matrix of the standard deviation: $\mathbf{D} = \text{diag}(\sigma_{\mathbf{X}^{\text{ho}}}, \sigma_{\mathbf{X}^{\text{hs}}}^{-1})$.

An observation of \mathbf{X}^{fa} is then given by $\mathbf{X}_t^{\text{ho}} + \mathbf{D} \cdot \mathbf{v}_{jk}$. To finalize, the empirical distribution $\hat{\mathbf{P}}_{\mathbf{X}^{\text{fa}}}$ can be calculated. Using
 515 this distribution, OTC can be applied to \mathbf{X}^{hs} and \mathbf{X}^{fa} to generate \mathbf{X}^{fa} . A more elaborate mathematical explanation can be
 found in Robin et al. (2019), a summary is given in Algorithm 4.

Algorithm 4 dOTC

Input:

Historical observations \mathbf{X}^{ho}
 Historical simulations \mathbf{X}^{hs}
 Future simulations \mathbf{X}^{fs}

Output:

Adjusted future simulations \mathbf{X}^{fa}

Calculate the empirical distributions $\hat{\mathbf{P}}_{\mathbf{X}^{\text{ho}}}$, $\hat{\mathbf{P}}_{\mathbf{X}^{\text{hs}}}$ and $\hat{\mathbf{P}}_{\mathbf{X}^{\text{fs}}}$

Calculate the optimal plan γ between $\hat{\mathbf{P}}_{\mathbf{X}^{\text{ho}}}$ and $\hat{\mathbf{P}}_{\mathbf{X}^{\text{hs}}}$

Calculate the optimal plan ϕ between $\hat{\mathbf{P}}_{\mathbf{X}^{\text{fs}}}$ and $\hat{\mathbf{P}}_{\mathbf{X}^{\text{hs}}}$

Calculate the Cholesky factor **D**

for all \mathbf{X}_t^{ho} **do**

Find cell \mathbf{c}_i containing \mathbf{X}_t^{ho}

Construct the vector of probabilities $\hat{\gamma}_{\mathbf{X}_t^{\text{ho}}} = (\gamma_{i,1}, \dots, \gamma_{i,J}) / p_{\mathbf{X}^{\text{ho}},i}$

Sample $j \in \{1, \dots, J\}$ according to the probability vector $\hat{\gamma}_{\mathbf{X}_t^{\text{ho}}}$

Construct the vector of probabilities $\hat{\phi}_{\mathbf{X}_t^{\text{hs}}} = (\phi_{j,1}, \dots, \phi_{j,K}) / p_{\mathbf{X}^{\text{hs}},j}$

Sample $k \in \{1, \dots, K\}$ according to the probability vector $\hat{\phi}_{\mathbf{X}_t^{\text{hs}}}$

Calculate the vector \mathbf{v}_{jk}

Calculate $\mathbf{X}_t^{\text{fa}} = \mathbf{X}_t^{\text{ho}} + \mathbf{D} \cdot \mathbf{v}_{jk}$

end for

Calculate the empirical distribution $\hat{\mathbf{P}}_{\mathbf{X}^{\text{fa}}}$

Apply OTC between \mathbf{X}^{fa} and \mathbf{X}^{fs} to generate \mathbf{X}^{fa}



3.4 Experimental design

Prior to all intensity-bias-adjusting methods, the thresholding occurrence-adjusting method was applied. In the intensity-bias-adjustment step, a balance was sought between randomness and computational power for the calculation of the intensity-bias-adjusting methods. Methods with randomised steps were repeated. As such, 10 calculations were made for dOTC. The resulting values of each index were averaged for further comparison. Biases on the indices were always calculated as raw or adjusted simulations minus observations, indicating a positive bias if the raw or adjusted simulations are larger than the observations and a negative bias if the simulations are smaller.

4 Results

In this section, the results will be shown first for the R index calculations for bias change, and then for the validation indices. For the validation indices, first the indices based on the adjusted variables are discussed, followed by an elaboration on the indices based on the derived variables. As the effect on discharge is the overarching goal of this paper and the discharge indices are affected by all other indices, those will be discussed last. All observed values and biases of both raw and adjusted simulations are presented in Table A1.

4.1 Bias change

The R index values for the variable averages, standard deviations and all indices are given in Table 2. The results vary considerably depending on the variable and/or index: for P, the bias can be considered almost stationary: only the 99.5th percentile has an R index above one. In contrast, for E, the R index values are above 1 for the middle percentiles and for the standard deviation, indicating some major changes in parts of the distribution, and consequently, the bias. For T, the mean and the lower extremes are clearly influenced, although the bias on the higher extremes does not change. The different effects on the variables are linked with the effect on the (cross-)correlations. For example, the lag-1 cross-correlation between P and E has an R index value of only 0.19, whereas the R index value for the cross-correlation between E and T is 1.20. Although the R index values are low for P, this does not imply that the R index values for the precipitation occurrence indices are low. With an R index value of 1.44, the auto-correlation bias clearly changes between both periods. However, this is not reflected by the other precipitation occurrence indices, which all but one have R index values lower than one.

Many of the R index values presented in Table 2 indicate that the bias changes between the two periods considered here (1970-1989 versus 1998-2017) might already be large enough to have an effect on the bias adjustment. As these periods are only separated by 10 years, this is an important indicator for the bias adjustment of late 21st century data, just as Chen et al. (2015) mentioned. However, it does not suffice to calculate just a few of these R index values. The results vary substantially among variables and even for the percentiles of a variable under consideration: while the 5th T percentile has an R index value of 2, the value for the 95th percentile is only 0.07. This could give an indication of why the methods perform more poorly for some of these indices. However, purely based on these results, it is impossible to say exactly what causes the



Table 2. R index values for 1970-1989 as historical period and 1998-2017 as future period

Precipitation		Temperature		Pot. Evaporation		Occurrence & Correlation	
Indices	R index	Indices	R index	Indices	R index	Indices	R index
P ₅	NaN	T ₅	2	E ₅	0.03	P _{lag1}	1.44
P ₂₅	0	T ₂₅	2	E ₂₅	0.47	P _{P00}	0.09
P ₅₀	0.10	T ₅₀	2	E ₅₀	1.47	P _{P10}	0.41
P ₇₅	0.13	T ₇₅	0.87	E ₇₅	2	N _{dry}	0.29
P ₉₀	0.19	T ₉₀	0.31	E ₉₀	1.14	corr _{E,T}	0.75
P ₉₅	0.17	T ₉₅	0.07	E ₉₅	1	corr _{P,E}	0.20
P ₉₉	0.58	T ₉₉	0.19	E ₉₉	0.47	corr _{P,T}	0
P _{99.5}	1.02	T _{99.5}	0.08	E _{99.5}	0.20	crosscorr _{E,T,0}	2
P _{mean}	0.18	T _{mean}	2	E _{mean}	1.06	crosscorr _{E,T,1}	0.90
P _{std}	0.72	T _{std}	0.50	E _{std}	2	crosscorr _{P,E,0}	0.31
						crosscorr _{P,E,1}	0.13
						crosscorr _{P,T,0}	0.10
						crosscorr _{P,T,1}	0.09

bias nonstationarities. Possible causes could be that recent trends such as those in precipitation extremes (Papalexiou and Montanari, 2019) are poorly captured by the models, that limiting mechanisms such as soil moisture (Bellprat et al., 2013) are poorly modelled or that natural variability influences (Addor and Fischer, 2015) the biases. However, discussing this in depth is out of the scope of the present study and deserves a separate study. In what follows, we will focus on the performance of the bias-adjusting methods and whether or not there is a link with these nonstationarities.

4.2 Precipitation amount

Figure 3 presents the RB_O and RB_{MB} values for the highest P percentiles. None of the residual bias values of the lower percentiles can be plotted as either the observations are 0 mm (P₅ and P₂₅) or the RB_O values are lower than zero (P₅₀). The percentiles could also have been plotted for wet days only (e.g. days with P higher than 0.1 mm/day), but as some methods change the number of dry days after the initial thresholding step, the dry days are also included in the calculation of the indices. This influences the RB_O and RB_{MB} values: they are generally higher when the dry days are not included.

The RB_O and RB_{MB} values depict a very similar performance for QDM, mQDM and MBCn, but a different performance for MRQNBC and dOTC. The similar performance of the former three is unsurprising, as their adjustments of P are all very similar. The only difference between QDM and MBCn versus mQDM is the time series to which the adjustment was applied, as the latter is based on the historical time series. QDM, mQDM and MBCn are consistently the best methods out of the five tested here, with the RB_{MB} values for P₇₅, P₉₀, P₉₅ and P₉₉ all below 0.5 and the RB_O values also below 0.5. The performances



of MRQNBC and dOTC are worse, but not poor either: P_{75} , P_{90} , P_{95} and P_{99} all have RB_O and RB_{MB} values lower than 1, but only for dOTC the majority of them (P_{75} , P_{90} and P_{95}) are below 0.5 for RB_{MB} . For both MRQNBC and dOTC, no RB_O values below 0.5 are obtained. As seen in Section 4.1, P was one of the few indices having almost all R index values below 1. This might be linked to the generally good results for P illustrated in this section.

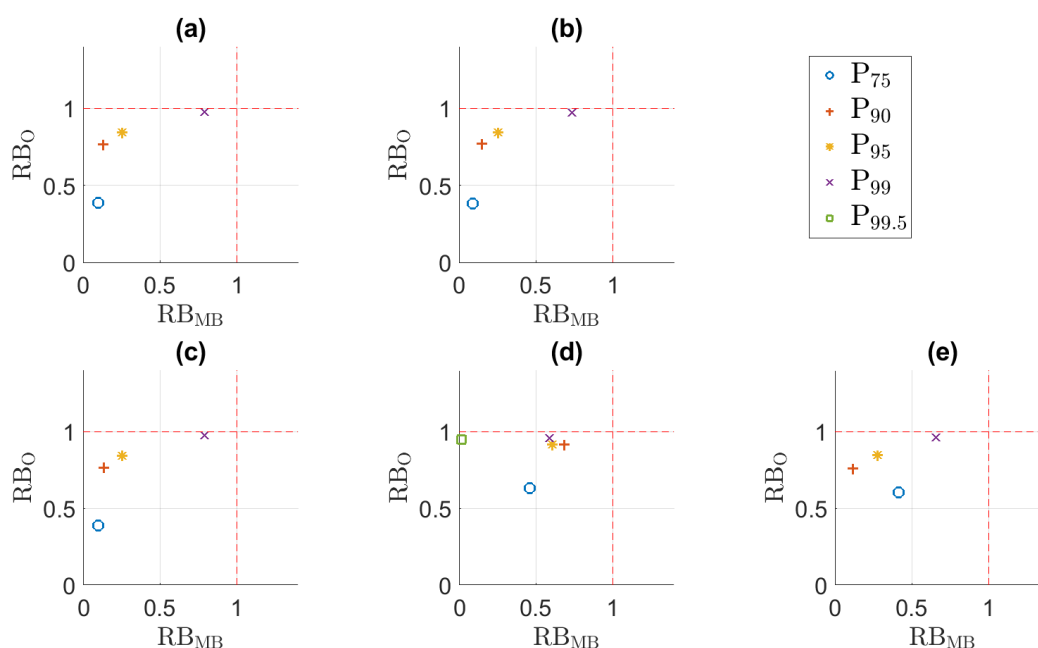


Figure 3. RB_{MB} versus RB_O for the precipitation indices. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC.

A surprising result for P is the high RB_{MB} value for $P_{99.5}$ for MRQNBC. This percentile is too biased for all other methods to be plotted. As the R index value was close to 1, it is possible that bias nonstationarity slightly influences the performance. For MRQNBC, however, the combination of QDM with the focus on correlation seemingly improves the performance of this percentile. As heavy precipitation values are clustered in time, the performance of the respective indices might be improved by the correlation. The good representation of heavy precipitation values in the MRQNBC-adjusted time series is also shown by P_{99} , for which MRQNBC has an RB_{MB} value of 0.59, the best of all methods.

4.3 Temperature

For the temperature adjustment, the RB_O and RB_{MB} values indicate that all methods result in a performance better than the raw climate simulations, except for MRQNBC (Fig. 4). In contrast to all other methods, only the residual bias values of T_{90} of MRQNBC are within the area delineated by the 1-1-lines. For all other indices, the bias is worse than in the climate simulations, with absolute biases up to 7°C . The results for MRQNBC are interesting, as T is the best understood variable (Shepherd, 2014)



and should thus not be hard to adjust. One line of reasoning could be the implementation of model trend preservation. While trend preservation is a prominent aspect in all other methods, the persistence preservation is at least as important in MRQNBC. The trade-off between both aspects of the bias adjustment thus seems to influence the result of MRQNBC, while the other methods can more easily adapt to and adjust the simulated T.

When comparing the indices for the other methods, the results are rather similar. They all have RB_O values close to 1, indicating that the bias difference is small in comparison to the absolute T values. Besides that, for every method, the lower T percentiles have RB_{MB} values that are too high to be plotted. However, despite their similar behaviour, the methods show some notable differences. The highest percentiles have the lowest RB_{MB} values for QDM and MBCn, which have the same percentiles by construction, but this differs for the other methods. For example, $T_{99.5}$ has an RB_{MB} value of 0.09 for QDM and MBCn, but only 0.43 for dOTC and 0.77 for mQDM. On the other hand, when considering all plotted percentiles, dOTC generally performs best. The highest RB_{MB} value for dOTC is 0.52 (T_{75}), whereas 0.65 (T_{75}) is the highest value for QDM and MBCn and 0.77 ($T_{99.5}$) for mQDM. Although broadly similar, the indices for QDM and mQDM display some interesting differences. Whereas for QDM T_{95} , T_{99} and $T_{99.5}$ have the lowest RB_{MB} values, T_{75} has the lowest value for mQDM. In contrast, QDM has the highest RB_{MB} value for T_{75} . This might imply that for the highest T values, it is better to follow the simulations, while for slightly lower values, it is better to only use the climate change signal. Yet, QDM has the best RB_{MB} values and might thus be preferable.

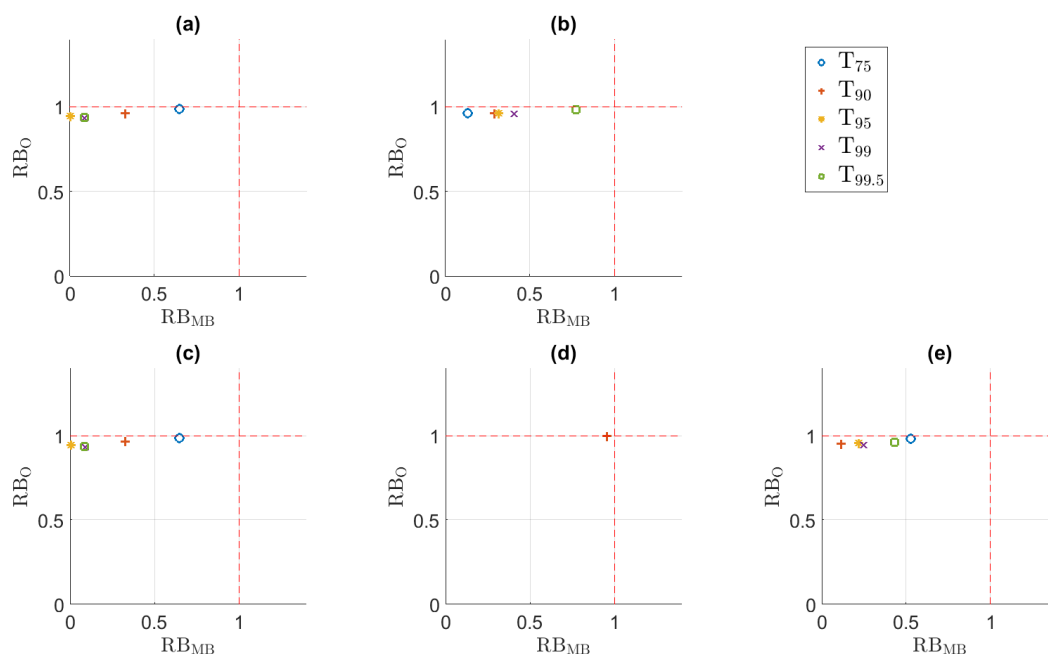


Figure 4. RB_{MB} versus RB_O for the temperature indices. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC.



595 For the lowest T values, all methods seem unable to handle the change in bias (as seen in Table 2): the RB_{MB} values are all higher than 1. This poor performance, combined with the high values for RB_O , might imply that it is better not to adjust T and work with the raw climate simulations. However, for the extreme T values, the absolute biases can be more than 1°C. Thus, depending on the research goal and the R index value, it might be important to consider whether or not T should be adjusted.

4.4 Potential evaporation

600 Figure 5 displays the RB_O and RB_{MB} values for the E indices. Only a few indices are shown for each method, or just one for dOTC, indicating that the performance after adjusting the bias is generally worse than the raw climate simulations. The indices plotted are E_{25} , E_{99} and $E_{99.5}$. The index E_5 also performs well, but cannot be plotted as its observed value is 0 mm. Thus, for the lowest and the highest percentiles, the bias-adjusting methods perform well, but they fail to capture the nonstationarity at the middle percentiles. These middle percentiles have high R index values: they are all greater than or equal to one. Only for
605 dOTC, it is possible to plot a percentile between E_{25} and E_{99} : E_{95} . However, for dOTC, this is also the only percentile for which it is possible to plot the RB_O and RB_{MB} values (respectively 1.00 and 0.86). For all other methods and all other percentiles, both RB_O and RB_{MB} values are higher than 1. The poor performance of dOTC might be related to its trend preservation characteristics. Of all methods, it is the one most explicitly designed to follow the simulated trend. This might thus imply that the nonstationarity for E is caused by poor model performance, although this should be investigated more in depth.

610 The percentiles that are plotted all have a high RB_O value, which is in this case caused by rather low biases to adjust. For example, $E_{99.5}$ had an observed value of 5.24 mm/day and only a bias of 0.27 mm/day, or 5%. There is consequently not much room for improvement, though the RB_{MB} values imply that the bias-adjusting method could be improved, except for $E_{99.5}$, which consistently has an RB_{MB} value lower than 0.5.

As in Section 4.3, there are interesting differences between QDM and mQDM. In contrast to the results for temperature,
615 mQDM performs better. For mQDM, the highest percentiles (E_{99} and $E_{99.5}$) have lower RB_{MB} values than for QDM. In this case, it thus seems better to only use the climate change signal. However, given the general poor performance of all methods, these results should be considered with care.

Given that the percentiles with a high R index value have a larger bias than the raw simulations after adjustment, and the added value for the other percentiles with respect to observed values is low, it can be advised not to adjust E. However, similar
620 to T, this should be evaluated on a case-by-case basis.

4.5 Correlation

When considering the correlation (Fig. 6), the methods generally perform well: most of the correlation indices can be plotted. However, there are some differences depending on the indices under consideration and the method. The indices that can always be plotted are the lag-0 cross-correlation between P and T, the lag-1 cross-correlation between P and T, the lag-1
625 cross-correlation between P and E and the correlation between P and T. Except for dOTC, all methods also perform well for the lag-0 cross-correlation between P and E. Yet, for the indices that can be plotted, the RB_O and RB_{MB} values show a considerable difference among the methods. For example, the lag-1 cross-correlation between P and E has RB_O values ranging

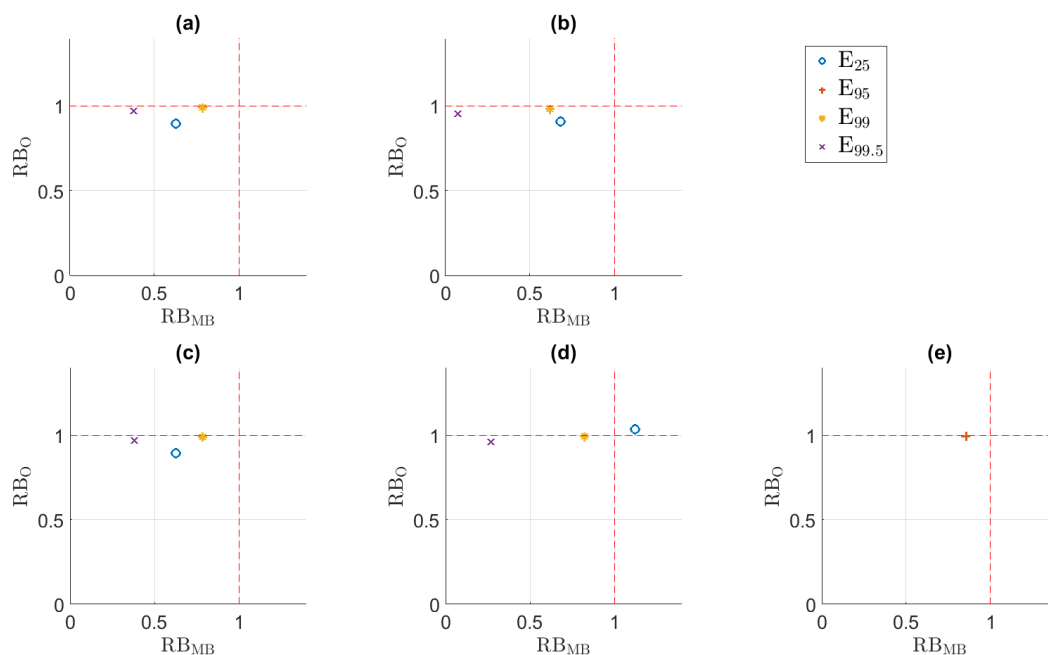


Figure 5. RB_{MB} versus RB_0 for the potential evaporation indices. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC.

from 0.29 (mQDM) to 0.69 (dOTC) and RB_{MB} values ranging from 0.08 (mQDM) to 0.60 (dOTC). The best method varies for each index: while dOTC does not perform well for most indices, it has the best performance for the correlation between P and T. It is interesting that the performance of dOTC seems to be either very good, or very poor. As dOTC is built around the idea of trend preservation and all-in-one adjustment, it is possible that the adjustment performs very well when the trend is properly modelled. Three out of the four correlations that dOTC adjusts well are based on T and P, the two variables that are well understood in the time frame under consideration. All indices that are not or less frequently present in the plots have one thing in common: they are based on the correlation between E and one of the other variables. The indices based on E thus consequently perform worst, except for the aforementioned lag-1 cross-correlation between P and E.

The correlation index performance seems to be related to the results of T (Section 4.3) and E (Section 4.4): correlations of E with another variable generally perform worse, and the (cross-)correlation between T and E with another variable performs the worst, in line with the R index values (Table 2). Although the multivariate bias-adjusting methods are supposed to adjust correlation, they seem to be unable to do so, as they generally have larger biases (though not on all indices) than the univariate bias-adjusting methods for the indices with the lowest residual bias values. This seems to indicate that the multivariate bias-adjusting methods, and especially MBCn and dOTC, are unable to adjust the correlation exactly because of the nonstationarity in the correlation that has to be overcome. In contrast, the univariate bias-adjusting methods neglect the adjustment of correlation and consequently do not have to overcome nonstationarity in the correlation bias. Yet, for the univariate bias-adjusting

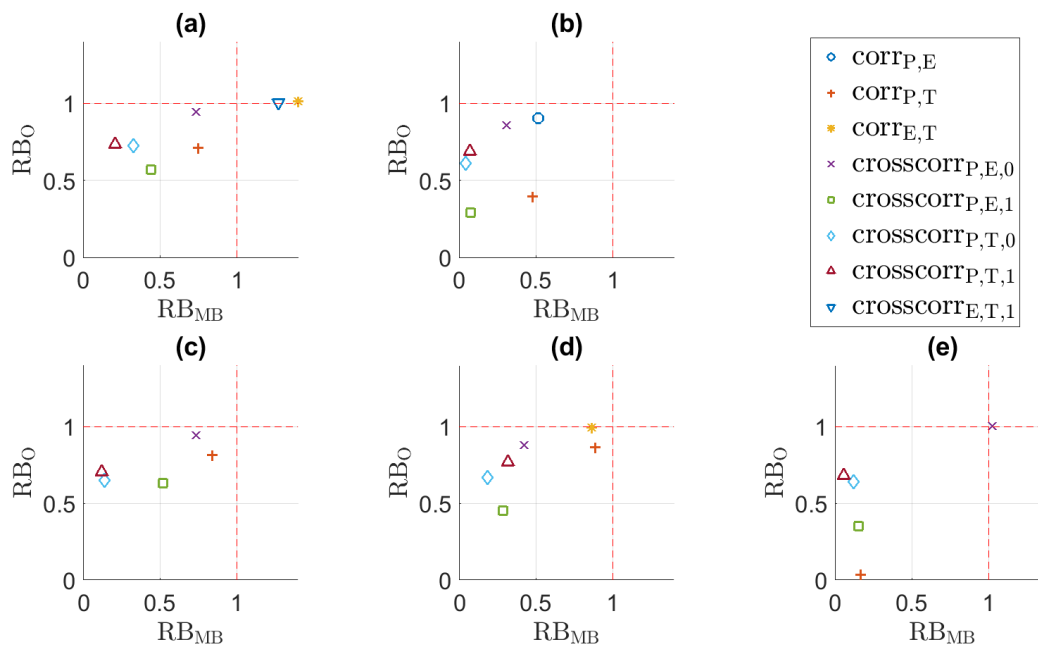


Figure 6. RB_{MB} versus RB_0 for the correlation indices. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC.

methods, the difference in adjustment of T and E seems to have an influence here as well, as illustrated by the different results
 645 for QDM and mQDM. Except for the correlations that have high RB_{MB} values for all methods, the results indicate that mQDM
 performs better. Thus, it might be better to use correlations of the observed time series than to adjust the simulated correlations.
 This is confirmed by the results for MRQNBC. Together with mQDM, this is the only method to have six indices with RB_{MB}
 values lower than one. Besides, it is also the only multivariate bias-adjusting method to have an RB_{MB} value for $corr_{E,T}$ lower
 than one, although it is only slightly lower.

650 4.6 Precipitation occurrence

Figure 7 displays the RB_0 and RB_{MB} values for the precipitation occurrence indices. When comparing these values, large
 differences among the methods and among the indices can be noted. The best-performing method seems to be QDM, as all
 the RB_{MB} values are lower than 0.5 and some RB_0 values are close to 0.5. When comparing the other methods, there is no
 clear difference between the univariate and multivariate bias-adjusting methods. The other univariate method, mQDM, and one
 655 multivariate method, MRQNBC, also perform better than the raw climate simulations for all indices. The other two methods,
 MBCn and dOTC, have respectively only one and two indices with both RB_0 and RB_{MB} values below 1.

Interestingly enough, the indices with RB_0 and RB_{MB} values below 1 are not the same for all methods. For the three best-
 performing methods, the dry-to-dry transition probability has very low RB_{MB} values (ranging from 0.3 to 0.18), while this

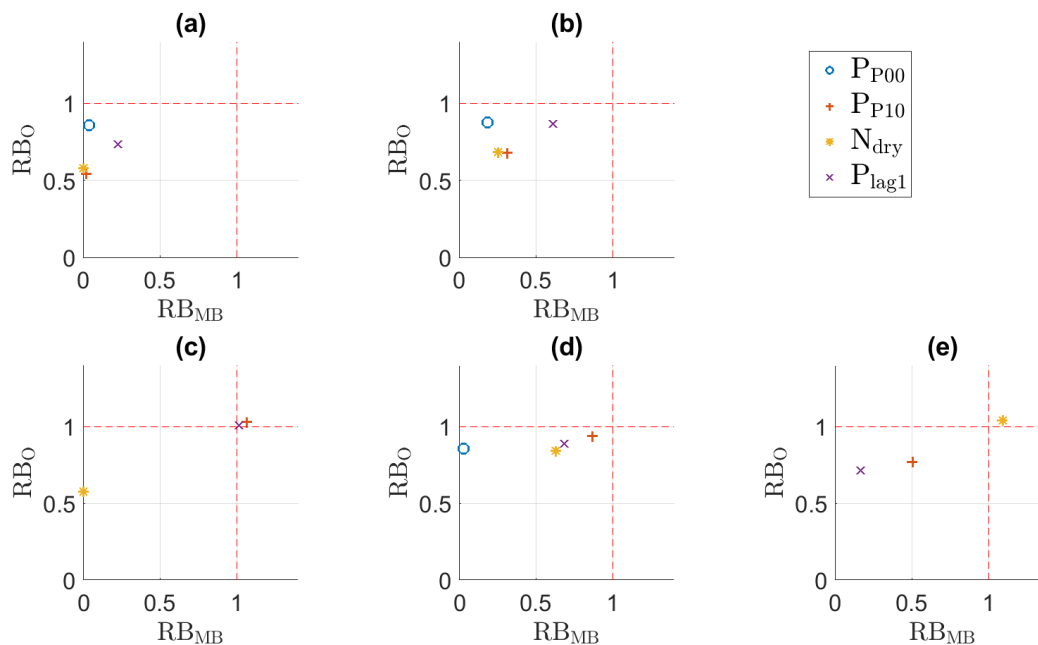


Figure 7. RB_{MB} versus RB_0 for the precipitation occurrence indices. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC.

index is absent from the MBCn and dOTC plots. The differences between those two plots are also notable. For MBCn, only the number of dry days has a very low RB_{MB} value (0), as the number of dry days is unaffected after the thresholding, whereas the lag-1 P auto-correlation and the wet-to-dry transition probability are more biased than the raw climate simulations. For dOTC, it is the other way around: the number of dry days is more biased after the application of dOTC, and the auto-correlation and the wet-to-dry transition probability perform well, with RB_{MB} values 0.50 or lower.

Another peculiar result that can be seen from Fig. 7 is the difference in dry day bias. Although all methods start with the same number of dry days, there are large differences among the RB_0 and RB_{MB} values for the number of dry days. The RB_0 values range from 0.58 to 1.04 and the RB_{MB} values from 0 to 1.09. QDM and MBCn perform best ($RB_{MB} = 0$), as the number of dry days is unaffected after the thresholding. For mQDM ($RB_{MB} = 0.25$), this holds by construction: instead of adjusting the threshold-adjusted climate model simulations, this method adjusts the observations. For MRQNBC ($RB_{MB} = 0.63$) and dOTC ($RB_{MB} = 1.09$), the results seem to imply that the multivariate framework of these methods has an influence on the number of dry days.

What the difference in transition probabilities implies for the time series, becomes more clear in Fig. 8. Although all adjusted simulations and the observations have more short wet spells than long ones, MBCn pronounces the short wet spell length more than the other methods, while the probability of longer wet spell lengths is lowered in comparison with other methods. Closest to the observations is mQDM. QDM and MRQNBC also perform well, a conclusion similar to that of Fig. 7.

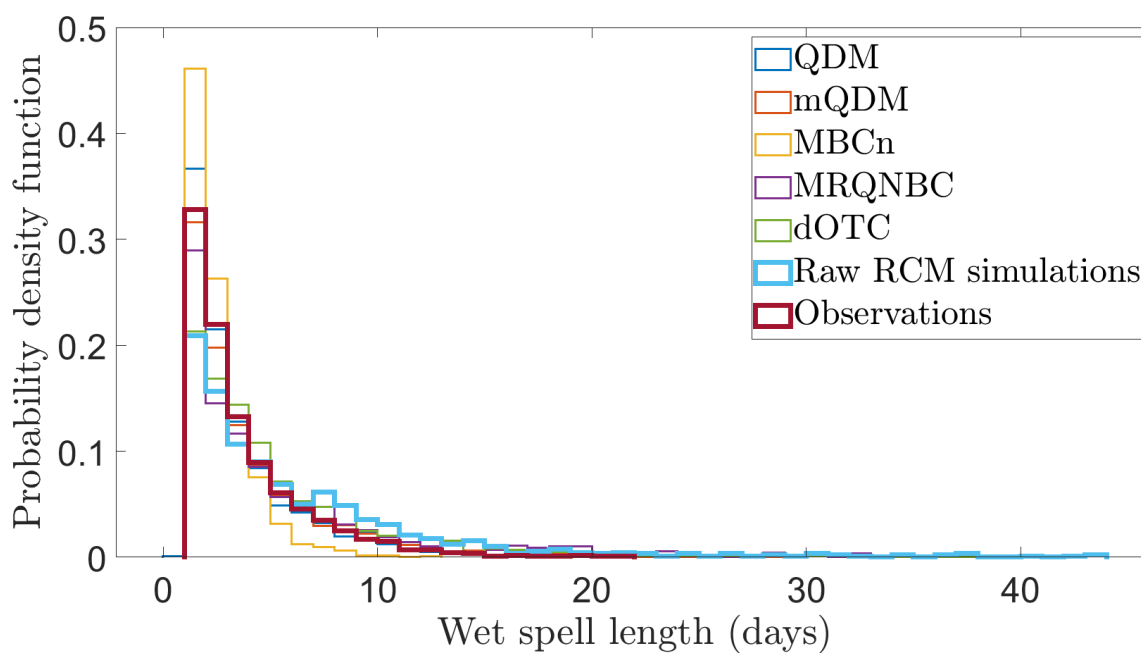


Figure 8. Wet spell length probability mass function for all adjusted simulations, the raw RCM simulations and the observations.

675 The difference in performance between QDM and the other methods seems to demonstrate that most strategies for retaining
a certain temporal structure or adjusting the temporal structure do not perform well. MRQNBC and mQDM depend heavily
on the temporal structure of the observations, and MBCn and dOTC have an important shuffling or recalculation aspect, all
of which lead to less reliable results at the end of the process. The poor performance of dOTC and MBCn for the temporal
structure was also discussed by François et al. (2020). As for mQDM and MRQNBC, it is notable that the temporal structure
680 does not change much from the calibration time series to the validation time series. At least, this is suggested by their relatively
good performance, which is based on using the observed time series (mQDM) and observed persistence statistics (MRQNBC).
Yet, this is no guarantee that these methods will be able to realistically adjust climate model simulations for the end of the
century.

Figure 7 also suggests that despite the high R index value, the P lag-1 auto-correlation is not necessarily poorly adjusted.
685 For QDM, this index has relatively low RB_O and RB_{MB} values. This could imply that the performance still depends on the
robustness to the bias nonstationarity of the methods under consideration. Or, as the other indices illustrate, the effect of bias
(non-)stationarity is not as large as the effect induced by the methods themselves. An example of this is the number of dry
days: though it has a low R index value, the performance varies substantially among the methods.



4.7 Discharge

690 All bias-adjusting methods perform better for the discharge percentiles compared to most other indices (Fig. 9). Although
the discharge is influenced by a combination of many effects, these appear to be small in the end result. For example, the
poor performance of E (Fig. 5) does not result in large discharge biases. Thus, it is the integration of precipitation amount,
precipitation occurrence and evaporation, and the routing effect that ultimately defines the resulting discharge. Generally,
many indices perform well, though there are some differences among the methods. The best-performing methods are QDM
695 and mQDM, as all indices have RB_{MB} values lower than 0.5 and some indices have RB_O values near to zero. For mQDM, the
20-year return period even has an RB_O value of -0.05. For MBCn, the results are also good. Most extreme values have RB_O and
 RB_{MB} values lower than 0.5; only the 5th percentile has an RB_O value of 0.96 and an RB_{MB} value of 0.89. As the only difference
between QDM and MBCn is the adjustment of occurrence, the results for discharge illustrates the importance of occurrence
adjustment. This variability in values seems to be a difference between the univariate and multivariate bias-adjusting methods.
700 For the two worst performing methods, i.e. MRQNBC and dOTC, some indices have RB_O and RB_{MB} values close to or higher
than 1 and some values between 0.5 and 1. These methods are thus unable to correctly adjust the bias for all indices. However,
although MRQNCB and dOTC seem to perform similarly, the indices with the worst RB_O and RB_{MB} values are different. For
MRQNBC, the 99.5th percentile and the 20-year return period have the highest values, whereas for dOTC, the 5th and 25th
percentile perform worse than the raw climate simulations. From the point of view of extreme discharges, dOTC is thus the
705 better method of these two. This might indicate that although not all occurrence indices of dOTC had lower RB_{MB} values than
those of MRQNBC, those that had (P_{lag1} and P_{P10}), had a larger influence on the extreme discharge values. Both indices are
partly based on the occurrence of wet days, and thus indicate that those need to be at the correct place in the time series for
extreme floods to be correctly simulated.

5 Discussion

710 In the previous section, the results for the bias adjustment by different methods and under climate change conditions were
reported. The effect of climate change on the bias was evaluated through the R index, which showed that the bias for some
indices cannot be considered stationary. For some of the indices (the lower percentiles of T and especially the middle percentiles
of E) the methods performed poorly, which could often be linked with the R index values. The methods clearly handle this
bias nonstationarity differently. It seems that the univariate bias-adjusting methods are far more robust: even for indices with
715 high R values, they are sometimes able to perform very well, with low RB_O and RB_{MB} values. This good performance thus
seems to imply that the more indices a bias-adjusting method directly adjusts, the more susceptible it is to problems related to
bias nonstationarity. However, this does not imply that QDM and mQDM are similar: while they are almost as good for many
variables, the poorer performance of mQDM for the precipitation occurrence indices is an indication that assuming that the
temporal structure of the past can be used for the future might be dangerous, as Johnson and Sharma (2011) and Kerkhoff et al.
720 (2014) already mentioned. Given that mQDM performed worse for two time periods separated by 10 years only, it is unlikely
that it is safe to use this method, or other delta change-based methods, for impact assessments targeting the end of the 21st

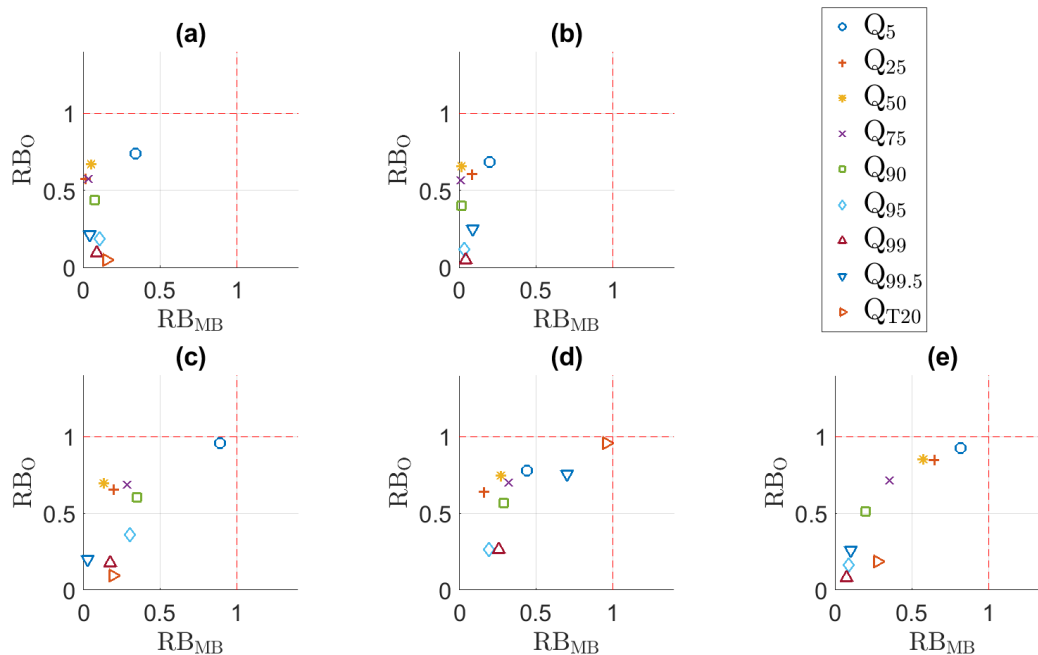


Figure 9. RB_{MB} versus RB_0 for the discharge percentiles and the 20-year return period value. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC.

century and depending on the temporal structure of time series. Yet, for some other indices, especially the correlation, mQDM performed better. Consequently, the exact choice should depend on the goals of the end user.

The results of the multivariate bias-adjusting methods too are not without nuance: though they are generally worse than the univariate bias-adjusting methods, their performance depends heavily on the variable under consideration and on the method itself. A clear example of this dependence on variables is the contrasting performance of dOTC to adjust T (Fig. 4) versus E (Fig. 9): the adjustment of E is much worse. This is a reminder that in a multivariate context, the multivariate methods are far less robust and can perform relatively good and poor at the same time for different variables. Therefore, there seems to be an interplay between the modelling of the variables and the method of calculation. Except for P, for which the results were similar, the methods performed differently for each variable. MRQNBC performed best in the context of temporal structure, for which it was designed (Fig. 7). For T, MBCn and dOTC performed better (Fig. 4). This could be related to their trend-preserving properties, which are more pronounced for those methods than for MRQNBC. For E (Fig. 5) and correlation (Fig. 6), dOTC displayed the most different results. For the former, the all-in-one and trend-preservation method did not seem robust enough. For the latter, it depended heavily on the type of correlation under consideration. These results seem to imply that the difference under bias-nonstationary conditions is not clear-cut for the different types of multivariate bias-adjusting methods. For the ‘marginal/dependence’ vs. the ‘all-in-one’ approach, consequently no clear conclusions can be drawn. For the amount of temporal alteration, it depends on the index under consideration. MRQNBC, which replaces the simulated correlations by



those of the observations performs well for the temporal structure, but performs worse for many other indices. For MBCn and dOTC, the effect of the difference in temporal alteration is less distinct and other properties, such as trend preservation, seem to have more influence.

To have a better view of how these results should be interpreted, the perspective of the end user should be considered (Maraun et al., 2015; Maraun and Widmann, 2018b). We used discharge as an example, using the relatively simple PDM. Although the residual bias values for all methods for E (Fig. 5) indicate a poor performance, the influence thereof on the discharge seems to be negligible (Fig. 9). Discharge is the variable that is of the highest importance for hydrological impact modelling, and the results indicate that most methods are able to adjust the forcing variables sufficiently in order to have a good simulation of discharge. However, the small differences between the methods should still be taken into account. Overall, QDM and mQDM perform best in adjusting the variables such that the discharge rates are the least biased in comparison with the observations. This is also important considering that bias adjustment can be applied for many different types of impact assessment. In other impact assessments, the differences could affect the result more than the discharge considered here. For example, forest fires (a typical compound event, discussed in a bias adjustment context in e.g. Yang et al. (2015), Cannon (2018), Zscheischler et al. (2019)) depend more heavily on T and E to simulate fire weather conditions. Besides such compound events, other applications are ecosystem functioning (Sippel et al., 2016), agriculture (Galmarini et al., 2019), or climate zone classification (Beck et al., 2018). In such studies the effect of bias nonstationarity can even be worse, whereas in other studies, depending more on P or other (so far) less-affected variables, the need for a bias nonstationarity-proof bias-adjusting method is less compelling. Anyway, the inability of some methods to adjust the biases in nonstationary conditions implies that a thorough assessment of possible bias nonstationarity should be made before bias-adjusting. If not done, the risk of reporting a wrong future projection is likely increased. Given the knowledge of bias nonstationarity, such uncertainties can be better characterised.

Returning to the discharge, it might be interesting to discuss whether or not the adjustment of E is truly needed. On the one hand, this variable is the most affected by bias nonstationarity. On the other hand, discharge is far less influenced by this variable than by P or temporal structure. The discharge has been calculated for this setting with raw E, the result of which is shown in Fig. 10. The results depend on the method: for QDM and mQDM, raw E data slightly exacerbate the results, while for dOTC the percentiles are all improved. Only for MRQNBC and MBCn, the results are highly dependent on the considered percentile. For MBCn, the 5th percentile and the 20-year return period value (with $RB_O \leq 0$) are improved, whereas the 95th and 99th percentile RB_O and RB_{MB} values are deteriorated. For MRQNBC, the results are opposite: the 5th percentile RB_O and RB_{MB} values are deteriorated, and the 95th and 99th values are improved.

The results for raw E seem to imply that, on the one hand and depending on the bias-adjusting method used, a well-considered choice of variables to adjust can give optimal results. On the other hand, the results demonstrate once more that the univariate methods are far more robust than the multivariate methods. Although the RB_O and RB_{MB} values are slightly deteriorated for QDM and mQDM in comparison with the discharge based on adjusted E, all values, and especially the values for the highest percentiles, still indicate a good bias adjustment. In general, an assessment like this can be done for the other types of impact studies discussed above, so that the influence of adjusting bias-nonstationary variables can be better understood.

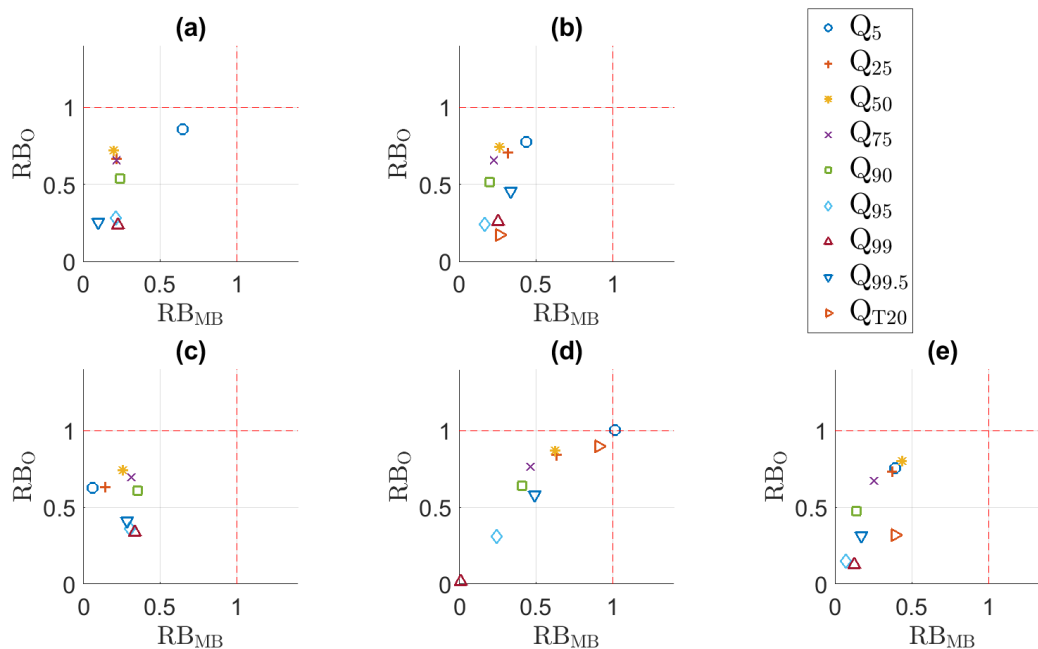


Figure 10. RB_{MB} versus RB_0 for the discharge percentiles and the 20-year return period value, calculated with raw evaporation. (a) QDM, (b) mQDM, (c) MBCn, (d) MRQNBC, (e) dOTC.

6 Conclusions

The goal of this paper was to assess how five bias-adjusting methods handle a climate change context with possible bias nonstationarity. Three of the methods were multivariate bias-adjusting methods: MRQNBC, MBCn and dOTC. The two other
 775 the bias-adjusting methods were univariate: one was a traditional bias-adjusting method, while the other was almost the same method, but modified according to the delta change paradigm. These univariate methods were used as a baseline to compare the multivariate bias-adjusting methods with. The climate change context, using 1970-1989 as calibration time period and 1998-2017 as validation time period, allowed us to calculate the change in bias between the periods, or the extent of bias stationarity, using the R index. All methods were calculated and compared using different indices, for which the residual biases relative to
 780 the observations and model bias were calculated.

The calculated R index values differed depending on the variable and variable index under consideration, but generally demonstrated that the bias of some of these indices is not stationary under climate change conditions. These changes could in some cases be clearly linked to the poor performance of bias-adjusting methods, such as for the lower percentiles of T or the middle percentiles of E. The performance was often poorer for the multivariate bias-adjusting methods, which corroborates the
 785 conclusions of Guo et al. (2020) that bias nonstationarity influences the performance of multivariate bias-adjusting methods. Although these methods have been developed during the last few years as a means to better adjust the biases, it seems that their



more complex calculations make them more vulnerable to bias nonstationarity. Thus, the univariate bias-adjusting methods, computationally less complex and not taking (potentially changing) correlations into account, seem to be more robust. Although effective difference in climate change impact is weakened by the hydrological model we used, the univariate methods still perform best. Studying other types of climate change impacts, the effect of bias nonstationarity could possible be even larger than discussed here.

The validation results could only be obtained by analysing and comparing a broad combination of indices. Considering only the mean or other standard statistics would have hidden many of the results seen. For example, in contrast to the results for the mean, the inclusion of both high and low extremes highlighted some problems with bias nonstationarity for some variables. As such, this study does not contradict earlier studies such as Maraun (2012), where the mean-based biases were found to be rather stable. Even a broader set of indices, such as the ETCCDI indices, was not enough to clearly discern between the methods. As such, we repeat the advice by Maraun and Widmann (2018a) to use indices not directly affected by bias-adjusting methods and to analyse the user needs before deciding upon the bias adjustment validation method. An important limitation is that we only used one GCM-RCM-combination. Using a model ensemble will be more informative, but could hide a single model's poor performance. On the other hand, similar assessments could also be used to discard poor-performing models (expanding upon methods such as those used in e.g. Brunner et al. (2019) or Tokarska et al. (2020)), based on the R index (also suggested by Maurer et al. (2013)) or the remaining bias after adjustment.

The results for the multivariate bias-adjusting methods assessed here are in line with François et al. (2020), especially for the problematic adjustment of occurrence. François et al. (2020) consequently state that the different multivariate bias-adjusting methods are based on different assumptions, and thus, the end user should make well-grounded choices on the method used. This also became clear in our assessment. However, François et al. (2020) did not study the effect of climate change and bias (non)stationarity and instead focused on model trend preservation, or trend nonstationarity. The results presented and discussed here, such as the contrasting results of MRQNBC and dOTC, imply that whether trend preservation was the focus of a method or not, can have an influence on the bias adjustment. However, it is yet unclear how trend nonstationarity and bias nonstationarity influence each other and how the most appropriate methods can be discerned, although it has been suggested to use trend-preserving methods whenever we can assume the models to correctly simulate the atmospheric processes (Maraun, 2016).

Although critical of their use, the results of this paper do not imply that multivariate bias-adjusting methods are not helpful. Many of the methods developed during the past few years can also be used for spatial bias adjustment, in which case the locations can be used as extra variables (see e.g. Vrac (2018)). A similar set-up has not been tested here, but the study by François et al. (2020) has proven the multivariate bias-adjusting methods to be very informative and robust for spatial adjustment: the spatial characteristics that influence local weather the most, such as orography.

Nonetheless, the results discussed in this paper indicate that many methods, and especially the multivariate bias-adjusting methods, fail in handling climate change and its resulting bias nonstationarity correctly. As authors have mentioned before (Ehret et al., 2012; Maraun, 2016; Nahar et al., 2017), this foremost implies that climate models have to become better at modelling the future: we need to be able to trust them as fully as possible. As long as this is not the case, bias adjustment



825 methods have to be developed that are more robust and that are able to help us assessing the future correctly. Yet, impact assessment cannot wait for new methods to be developed and/or tested: we need to prepare ourselves for the future as soon as possible. As was shown here, we can state for the current generation of methods that the fewer assumptions and calculations a method needs, the more robust it is when used in a climate change context. Given this statement, we advise to use univariate bias-adjusting methods, until it becomes more clear how it can be ensured that multivariate methods certainly perform well in a climate change context.

830 *Code and data availability.* The code for the computations is publicly available at <https://doi.org/10.5281/zenodo.4247518> (Hydro-Climate Extremes Lab – Ghent University, 2020). The RCA4 data are downloaded and are available from the Earth System Grid Federation data repository. The local observations were obtained from RMI in Belgium, and cannot be shared with third parties.

Author contributions. JVDV, BDB and NV designed the experiments. JVDV developed the code and performed the calculations. JVDV prepared the manuscript with contributions from MD, BDB and NV. All co-authors contributed to the interpretation of the results.

Competing interests. The authors declare that they have no conflict of interests.

835 *Acknowledgements.* J. Van de Velde would like to thank Y. Robin and R. Mehrotra for some helpful discussion on the use of respectively dOTC and MRQNBC. The authors are grateful to the RMI for allowing the use of 117-year Uccle dataset. This work was funded by FWO, grant number G.0039.18N.



Appendix A: Appendix A

Table A1: Observed values, and biases for the raw and adjusted climate simulations.

Index	Observed value	Bias					
		Raw climate simulations	QDM	mQDM	MBCn	MRQNBC	dOTC
P ₅ (mm)	0.00	0.00	0.00	0.00	0.00	-0.40	0
P ₂₅ (mm)	0.00	0.08	0.00	0.00	0.00	0.00	0.42
P ₅₀ (mm)	0.10	1.01	0.05	0.10	0.05	0.27	0.83
P ₇₅ (mm)	2.70	1.83	-0.18	-0.17	-0.18	-0.84	0.76
P ₉₀ (mm)	7.40	1.99	-0.26	-0.30	-0.26	-1.36	-0.23
P ₉₅ (mm)	11.42	2.38	-0.61	-0.60	-0.61	-1.44	-0.65
P ₉₉ (mm)	21.80	2.36	-1.86	-1.73	-1.86	1.38	-1.55
P _{99.5} (mm)	29.09	1.56	-4.20	-3.97	-4.20	0.02	-4.02
T ₅ (°C)	0.40	-0.31	-0.70	-1.43	-0.70	-0.63	-0.94
T ₂₅ (°C)	6.30	-0.08	-0.68	-1.55	-0.68	-3.00	-0.73
T ₅₀ (°C)	11.40	-0.40	-0.81	-0.88	-0.81	-3.24	-0.70
T ₇₅ (°C)	16.10	-0.70	-0.46	0.09	-0.46	-1.07	-0.37
T ₉₀ (°C)	19.40	-1.07	-0.35	0.31	-0.35	1.02	0.12
T ₉₅ (°C)	21.30	-1.17	-0.01	0.37	-0.01	2.59	0.25
T ₉₉ (°C)	24.95	-1.85	-0.16	-0.75	-0.16	5.95	0.46
T _{99.5} (°C)	25.90	-1.80	0.16	-1.39	0.16	7.01	0.77
E ₅ (mm)	0.00	0.20	0.00	0.00	0.00	-0.04	0
E ₂₅ (mm)	0.52	0.15	-0.09	-0.10	-0.09	-0.16	-0.52
E ₅₀ (mm)	1.42	0.05	-0.27	-0.28	-0.27	-0.38	-0.77
E ₇₅ (mm)	2.69	-0.02	-0.34	-0.35	-0.34	-0.58	-0.65
E ₉₀ (mm)	3.65	0.10	-0.27	-0.27	-0.27	-0.47	-0.18
E ₉₅ (mm)	4.21	0.15	-0.30	-0.28	-0.30	-0.41	0.125
E ₉₉ (mm)	5.02	0.21	-0.16	-0.13	-0.16	-0.17	0.69
E _{99.5} (mm)	5.24	0.27	-0.10	-0.02	-0.10	-0.07	1.03
corr _{P,E} (-)	-0.18	-0.04	-0.06	0.02	0.19	0.17	0.57
corr _{P,T} (-)	-0.16	0.18	0.14	0.09	0.16	0.16	0.04
corr _{E,T} (-)	0.82	-0.02	-0.03	-0.09	-0.84	0.02	-0.45



Table A1: (continued)

$\text{crosscorr}_{P,E,0}(-)$	0.30	0.06	-0.04	-0.02	0.05	-0.03	0.12
$\text{crosscorr}_{P,E,1}(-)$	0.24	0.19	0.08	-0.01	0.10	0.05	0.11
$\text{crosscorr}_{P,T,0}(-)$	0.36	0.15	0.05	0.01	0.02	-0.03	0.06
$\text{crosscorr}_{P,T,1}(-)$	0.38	0.13	0.02	-0.01	0.02	-0.04	0.08
$\text{crosscorr}_{E,T,0}(-)$	0.93	-0.00	-0.02	-0.05	-0.29	-0.02	-0.21
$\text{crosscorr}_{E,T,1}(-)$	0.91	0.01	-0.01	-0.04	-0.27	-0.01	-0.24
$P_{P00}(-)$	0.65	-0.10	-0.00	-0.02	-0.17	0.00	-0.37
$P_{P10}(-)$	0.32	-0.15	0.00	-0.05	0.16	-0.13	-0.07
$N_{\text{dry}}(-)$	3470.00	-1466.00	0.00	-373.00	0.00	-923.00	-1604.20
$P_{\text{lag1}}(-)$	0.33	0.11	0.02	0.07	-0.12	0.08	0.05
$Q_5(\text{m}^3/\text{s})$	2.30	0.92	-0.32	-0.18	0.82	-0.40	1.50
$Q_{25}(\text{m}^3/\text{s})$	3.36	1.45	0.02	-0.12	0.29	-0.23	1.50
$Q_{50}(\text{m}^3/\text{s})$	4.39	1.53	0.08	0.02	-0.20	-0.42	1.33
$Q_{75}(\text{m}^3/\text{s})$	5.72	2.52	-0.08	-0.03	-0.72	-0.81	1.51
$Q_{90}(\text{m}^3/\text{s})$	7.83	4.76	-0.36	-0.07	-1.66	-1.37	2.12
$Q_{95}(\text{m}^3/\text{s})$	10.09	9.22	-1.00	-0.33	-2.78	-1.78	2.94
$Q_{99}(\text{m}^3/\text{s})$	18.71	18.58	-1.65	-0.78	-3.21	4.77	5.77
$Q_{99.5}(\text{m}^3/\text{s})$	23.90	19.70	0.84	-1.77	-0.57	13.81	6.61
$Q_{T20}(\text{m}^3/\text{s})$	48.69	54.61	8.36	-3.41	-10.40	52.45	25.03



References

- Addor, N. and Fischer, E. M.: The influence of natural variability and interpolation errors on bias characterization in RCM simulations, *Journal of Geophysical Research: Atmospheres*, 120, 10–180, <https://doi.org/10.1002/2014JD022824>, 2015.
- Addor, N. and Seibert, J.: Bias correction for hydrological impact studies – beyond the daily perspective, *Hydrological Processes*, 28, 4823–4828, <https://doi.org/10.1002/hyp.10238>, 2014.
- Argüeso, D., Evans, J. P., and Fita, L.: Precipitation bias correction of very high resolution regional climate models, *Hydrology and Earth System Sciences*, 17, 4379, <https://doi.org/10.5194/hess-17-4379-2013>, 2013.
- 845 Bárdossy, A. and Pegram, G.: Multiscale spatial recorrelation of RCM precipitation to produce unbiased climate change scenarios over large areas and small, *Water Resources Research*, 48, W09 502, <https://doi.org/10.1029/2011WR011524>, 2012.
- Beck, H. E., Zimmermann, N. E., McVicar, T. R., Vergopolan, N., Berg, A., and Wood, E. F.: Present and future Köppen-Geiger climate classification maps at 1-km resolution, *Scientific data*, 5, 180 214, <https://doi.org/10.1038/sdata.2018.214>, 2018.
- Bellprat, O., Kotlarski, S., Lüthi, D., and Schär, C.: Physical constraints for temperature biases in climate models, *Geophysical Research*
850 *Letters*, 40, 4042–4047, <https://doi.org/10.1002/grl.50737>, 2013.
- Berg, P., Feldmann, H., and Panitz, H.-J.: Bias correction of high resolution regional climate model data, *Journal of Hydrology*, 448, 80–92, <https://doi.org/10.1016/j.jhydrol.2012.04.026>, 2012.
- Brunner, L., Lorenz, R., Zumwald, M., and Knutti, R.: Quantifying uncertainty in European climate projections using combined performance-independence weighting, *Environmental Research Letters*, 14, 124 010, <https://doi.org/10.1088/1748-9326/ab492f>, 2019.
- 855 Buser, C. M., Künsch, H. R., Lüthi, D., Wild, M., and Schär, C.: Bayesian multi-model projection of climate: bias assumptions and interannual variability, *Climate Dynamics*, 33, 849–868, <https://doi.org/10.1007/s00382-009-0588-6>, 2009.
- Cabus, P.: River flow prediction through rainfall–runoff modelling with a probability-distributed model (PDM) in Flanders, Belgium, *Agricultural Water Management*, 95, 859–868, <https://doi.org/10.1016/j.agwat.2008.02.013>, 2008.
- Cannon, A. J.: Multivariate bias correction of climate model output: Matching marginal distributions and intervariable dependence structure,
860 *Journal of Climate*, 29, 7045–7064, <https://doi.org/10.1175/JCLI-D-15-0679.1>, 2016.
- Cannon, A. J.: Multivariate quantile mapping bias correction: an N-dimensional probability density function transform for climate model simulations of multiple variables, *Climate Dynamics*, 50, 31–49, <https://doi.org/10.1007/s00382-017-3580-6>, 2018.
- Cannon, A. J., Sobie, S. R., and Murdock, T. Q.: Bias correction of GCM precipitation by quantile mapping: How well do methods preserve changes in quantiles and extremes?, *Journal of Climate*, 28, 6938–6959, <https://doi.org/10.1175/JCLI-D-14-00754.1>, 2015.
- 865 Chen, J., Brissette, F. P., and Lucas-Picher, P.: Assessing the limits of bias-correcting climate model outputs for climate change impact studies, *Journal of Geophysical Research: Atmospheres*, 120, 1123–1136, <https://doi.org/10.1002/2014JD022635>, 2015.
- Christensen, J. H., Boberg, F., Christensen, O. B., and Lucas-Picher, P.: On the need for bias correction of regional climate change projections of temperature and precipitation, *Geophysical Research Letters*, 35, L20 709, <https://doi.org/10.1029/2008GL035694>, 2008.
- Clark, M., Gangopadhyay, S., Hay, L., Rajagopalan, B., and Wilby, R.: The Schaake shuffle: A method for reconstructing space–time
870 variability in forecasted precipitation and temperature fields, *Journal of Hydrometeorology*, 5, 243–262, [https://doi.org/10.1175/1525-7541\(2004\)005<0243:TSSAMF>2.0.CO;2](https://doi.org/10.1175/1525-7541(2004)005<0243:TSSAMF>2.0.CO;2), 2004.
- Cuturi, M.: Sinkhorn distances: Lightspeed computation of optimal transport, in: *Advances in Neural Information Processing Systems*, pp. 2292–2300, 2013.



- De Jongh, I. L. M., Verhoest, N. E. C., and De Troch, F. P.: Analysis of a 105-year time series of precipitation observed at Uccle, Belgium,
875 International Journal of Climatology, 26, 2023–2039, <https://doi.org/10.1002/joc.1352>, 2006.
- Dekens, L., Parey, S., Grandjacques, M., and Dacunha-Castelle, D.: Multivariate distribution correction of climate model outputs: A generalization of quantile mapping approaches, *Environmetrics*, 28, e2454, <https://doi.org/10.1002/env.2454>, 2017.
- Demarée, G. R.: The centennial recording raingauge of the Uccle Plateau: Its history, its data and its applications, *Houille Blanche*, 4, 95–102, 2003.
- 880 Derbyshire, J.: The siren call of probability: Dangers associated with using probability for consideration of the future, *Futures*, 88, 43–54, <https://doi.org/10.1016/j.futures.2017.03.011>, 2017.
- Di Luca, A., de Elía, R., and Laprise, R.: Challenges in the quest for added value of regional climate dynamical downscaling, *Current Climate Change Reports*, 1, 10–21, <https://doi.org/10.1007/s40641-015-0003-9>, 2015.
- Eberhart, R. and Kennedy, J.: A new optimizer using Particle Swarm Theory, in: *Proceedings of the Sixth International Symposium on Micro*
885 *Machine and Human Science*, pp. 39–43, IEEE, 1995.
- Ehret, U., Zehe, E., Wulfmeyer, V., Warrach-Sagi, K., and Liebert, J.: HESS Opinions" Should we apply bias correction to global and regional climate model data?", *Hydrology and Earth System Sciences*, 16, 3391–3404, <https://doi.org/10.5194/hess-16-3391-2012>, 2012.
- Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., and Taylor, K. E.: Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental design and organization, *Geoscientific Model Development*, 9, 1937–1958,
890 <https://doi.org/10.5194/gmd-9-1937-2016>, 2016.
- Fosser, G., Kendon, E. J., Stephenson, D., and Tucker, S.: Convection-permitting models offer promise of more certain extreme rainfall projections, *Geophysical Research Letters*, p. e2020GL088151, <https://doi.org/10.1029/2020GL088151>, 2020.
- François, B., Vrac, M., Cannon, A. J., Robin, Y., and Allard, D.: Multivariate bias corrections of climate simulations: Which benefits for which losses?, *Earth System Dynamics*, 2020, 1–41, <https://doi.org/10.5194/esd-11-537-2020>, 2020.
- 895 Galmarini, S., Cannon, A., Ceglar, A., Christensen, O., Dentener, F. J., de Noblet-Ducoudré, N., Doblas-Reyes, F. J., and Vrac, M.: Adjusting climate model bias for agricultural impact assessment: How to cut the mustard, *Climate Services*, 13, 65–69, <https://doi.org/10.1016/j.cliser.2019.01.004>, 2019.
- Gudmundsson, L., Bremnes, J. B., Haugen, J. E., and Engen-Skaugen, T.: Downscaling RCM precipitation to the station scale using statistical transformations—a comparison of methods, *Hydrology and Earth System Sciences*, 16, 3383–3390, [https://doi.org/10.5194/hess-16-3383-](https://doi.org/10.5194/hess-16-3383-2012)
900 [2012](https://doi.org/10.5194/hess-16-3383-2012), 2012.
- Guo, Q., Chen, J., Zhang, X. J., Xu, C.-Y., and Chen, H.: Impacts of using state-of-the-art multivariate bias correction methods on hydrological modeling over North America, *Water Resources Research*, 56, e2019WR026659, <https://doi.org/10.1029/2019WR026659>, 2020.
- Gutiérrez, J. M., Maraun, D., Widmann, M., Huth, R., Hertig, E., Benestad, R., Roessler, O., Wibig, J., Wilcke, R., Kotlarski, S., et al.: An intercomparison of a large ensemble of statistical downscaling methods over Europe: results from the VALUE perfect predictor cross-validation experiment, *International Journal of Climatology*, 39, 3750–3785, <https://doi.org/10.1002/joc.5462>, 2019.
- 905 Gutjahr, O. and Heinemann, G.: Comparing precipitation bias correction methods for high-resolution regional climate simulations using COSMO-CLM, *Theoretical and Applied Climatology*, 114, 511–529, <https://doi.org/10.1007/s00704-013-0834-z>, 2013.
- Gutowski, W. J., Decker, S. G., Donavon, R. A., Pan, Z., Arritt, R. W., and Takle, E. S.: Temporal–spatial scales of observed and simulated precipitation in central US climate, *Journal of Climate*, 16, 3841–3847, [https://doi.org/10.1175/1520-](https://doi.org/10.1175/1520-0442(2003)016<3841:TSOAS>2.0.CO;2)
910 [0442\(2003\)016<3841:TSOAS>2.0.CO;2](https://doi.org/10.1175/1520-0442(2003)016<3841:TSOAS>2.0.CO;2), 2003.



- Haerter, J., Hagemann, S., Moseley, C., and Piani, C.: Climate model bias correction and the role of timescales, *Hydrology and Earth System Sciences*, 15, 1065–1073, <https://doi.org/10.5194/hess-15-1065-2011>, 2011.
- Hagemann, S., Chen, C., Haerter, J. O., Heinke, J., Gerten, D., and Piani, C.: Impact of a statistical bias correction on the projected hydrological changes obtained from three GCMs and two hydrology models, *Journal of Hydrometeorology*, 12, 556–578, <https://doi.org/10.1175/2011JHM1336.1>, 2011.
- 915
- Hakala, K., Addor, N., and Seibert, J.: Hydrological modeling to evaluate climate model simulations and their bias correction, *Journal of Hydrometeorology*, 19, 1321–1337, <https://doi.org/10.1175/JHM-D-17-0189.1>, 2018.
- Hay, L. E. and Clark, M. P.: Use of statistically and dynamically downscaled atmospheric model output for hydrologic simulations in three mountainous basins in the western United States, *Journal of Hydrology*, 282, 56–75, [https://doi.org/10.1016/S0022-1694\(03\)00252-X](https://doi.org/10.1016/S0022-1694(03)00252-X),
- 920 2003.
- Helsen, S., van Lipzig, N. P. M., Demuzere, M., Vanden Broucke, S., Caluwaerts, S., De Cruz, L., De Troch, R., Hamdi, R., Termonia, P., Van Schaeybroeck, B., and Wouters, H.: Consistent scale-dependency of future increases in hourly extreme precipitation in two convection-permitting climate models, *Climate Dynamics*, 54, 1–14, <https://doi.org/10.1007/s00382-019-05056-w>, 2019.
- Hewitson, B. C., Daron, J., Crane, R. G., Zermoglio, M. F., and Jack, C.: Interrogating empirical-statistical downscaling, *Climatic change*,
- 925 122, 539–554, <https://doi.org/10.1007/s10584-013-1021-z>, 2014.
- Higham, N. J.: Computing a nearest symmetric positive semidefinite matrix, *Linear Algebra and its Applications*, 103, 103–118, [https://doi.org/10.1016/0024-3795\(88\)90223-6](https://doi.org/10.1016/0024-3795(88)90223-6), 1988.
- Ho, C. K., Stephenson, D. B., Collins, M., Ferro, C. A. T., and Brown, S. J.: Calibration strategies: a source of additional uncertainty in climate change projections, *Bulletin of the American Meteorological Society*, 93, 21–26, <https://doi.org/10.1175/2011BAMS3110.1>, 2012.
- 930 Hui, Y., Chen, J., Xu, C.-Y., Xiong, L., and Chen, H.: Bias nonstationarity of global climate model outputs: The role of internal climate variability and climate model sensitivity, *International Journal of Climatology*, 39, 2278–2294, <https://doi.org/10.1002/joc.5950>, 2019.
- Hydro-Climate Extremes Lab – Ghent University: h-cel/ImpactofBiasNonstationarity: Impact of bias nonstationarity: calculations, <https://doi.org/10.5281/zenodo.4247518>, 2020.
- Ines, A. V. M. and Hansen, J. W.: Bias correction of daily GCM rainfall for crop simulation studies, *Agricultural and Forest Meteorology*,
- 935 138, 44–53, <https://doi.org/10.1016/j.agrformet.2006.03.009>, 2006.
- IPCC: Managing the Risks of Extreme Events and Disasters to Advance Climate Change Adaptation: Special Report of the Intergovernmental Panel on Climate Change, Cambridge University Press, 2012.
- IPCC: Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change, Cambridge University Press, 2013.
- 940 Jacob, D., Petersen, J., Eggert, B., Alias, A., Christensen, O. B., Bouwer, L. M., Braun, A., Colette, A., Déqué, M., Georgievski, G., Georgopoulou, E., Gobiet, A., Menut, L., Nikulin, G., Haensler, A., Hempelmann, N., Jones, C., Keuler, K., Kovats, S., Kröner, N., Kotlarski, S., Kriegsmann, A., Martin, E., van Meijgaard, E., Moseley, C., Pfeifer, S., Preuschmann, S., Radermacher, C., Radtke, K., Rechid, D., Rounsevell, M., Samuelsson, P., Somot, S., Soussana, J.-F., Teichmann, C., Valentini, R., Vautard, R., Weber, B., and Yiou, P.: EURO-CORDEX: new high-resolution climate change projections for European impact research, *Regional Environmental Change*, 14, 563–578,
- 945 <https://doi.org/10.1007/s10113-013-0499-2>, 2014.
- Johnson, F. and Sharma, A.: Accounting for interannual variability: A comparison of options for water resources climate change impact assessments, *Water Resources Research*, 47, W04 508, <https://doi.org/10.1029/2010WR009272>, 2011.



- Johnson, F. and Sharma, A.: A nesting model for bias correction of variability at multiple time scales in general circulation model precipitation simulations, *Water Resources Research*, 48, W01 504, <https://doi.org/10.1029/2011WR010464>, 2012.
- 950 Kendon, E. J., Ban, N., Roberts, N. M., Fowler, H. J., Roberts, M. J., Chan, S. C., Evans, J. P., Fosse, G., and Wilkinson, J. M.: Do convection-permitting regional climate models improve projections of future precipitation change?, *Bulletin of the American Meteorological Society*, 98, 79–93, <https://doi.org/10.1175/BAMS-D-15-0004.1>, 2017.
- Kerkhoff, C., Künsch, H. R., and Schär, C.: Assessment of bias assumptions for climate models, *Journal of Climate*, 27, 6799–6818, <https://doi.org/10.1175/JCLI-D-13-00716.1>, 2014.
- 955 Klemeš, V.: Operational testing of hydrological simulation models, *Hydrological Sciences Journal*, 31, 13–24, <https://doi.org/10.1080/02626668609491024>, 1986.
- Knol, D. L. and ten Berge, J. M. F.: Least-squares approximation of an improper correlation matrix by a proper one, *Psychometrika*, 54, 53–61, <https://doi.org/10.1007/BF02294448>, 1989.
- Kotlarski, S., Keuler, K., Christensen, O. B., Colette, A., Déqué, M., Gobiet, A., Goergen, K., Jacob, D., Lüthi, D., Van Meijgaard, E., et al.:
960 Regional climate modeling on European scales: a joint standard evaluation of the EURO-CORDEX RCM ensemble, *Geoscientific Model Development*, 7, 1297–1333, <https://doi.org/10.5194/gmd-7-1297-2014>, 2014.
- Lenderink, G., Buishand, A., and Van Deursen, W.: Estimates of future discharges of the river Rhine using two scenario methodologies: direct versus delta approach, *Hydrology and Earth System Sciences*, 11, 1145–1159, <https://doi.org/10.5194/hess-11-1145-2007>, 2007.
- Li, C., Sinha, E., Horton, D. E., Diffenbaugh, N. S., and Michalak, A. M.: Joint bias correction of temperature and precipitation in climate
965 model simulations, *Journal of Geophysical Research: Atmospheres*, 119, 13–153, <https://doi.org/10.1002/2014JD022514>, 2014.
- Li, H., Sheffield, J., and Wood, E. F.: Bias correction of monthly precipitation and temperature fields from Intergovernmental Panel on Climate Change AR4 models using equidistant quantile matching, *Journal of Geophysical Research: Atmospheres*, 115, D10 101, <https://doi.org/10.1029/2009JD012882>, 2010.
- Maraun, D.: Nonstationarities of regional climate model biases in European seasonal mean temperature and precipitation sums, *Geophysical
970 Research Letters*, 39, <https://doi.org/10.1029/2012GL051210>, 2012.
- Maraun, D.: Bias correcting climate change simulations—a critical review, *Current Climate Change Reports*, 2, 211–220, <https://doi.org/10.1007/s40641-016-0050-x>, 2016.
- Maraun, D. and Widmann, M.: Cross-validation of bias-corrected climate simulations is misleading, *Hydrology and Earth System Sciences*, 22, 4867–4873, <https://doi.org/10.5194/hess-22-4867-2018>, 2018a.
- 975 Maraun, D. and Widmann, M.: Statistical Downscaling and Bias Correction for Climate Research, Cambridge University Press, <https://doi.org/10.1017/9781107588783>, 2018b.
- Maraun, D., Widmann, M., Gutiérrez, J. M., Kotlarski, S., Chandler, R. E., Hertig, E., Wibig, J., Huth, R., and Wilcke, R. A. I.: VALUE: A framework to validate downscaling approaches for climate change studies, *Earth’s Future*, 3, 1–14, <https://doi.org/10.1002/2014EF000259>, 2015.
- 980 Matalas, N. C.: Mathematical assessment of synthetic hydrology, *Water Resources Research*, 3, 937–945, <https://doi.org/10.1029/WR003i004p00937>, 1967.
- Maurer, E. P., Das, T., and Cayan, D. R.: Errors in climate model daily precipitation and temperature output: time invariance and implications for bias correction, *Hydrology and Earth System Sciences*, 17, 2147–2159, <https://doi.org/10.5194/hess-17-2147-2013>, 2013.
- Mehrotra, R. and Sharma, A.: An improved standardization procedure to remove systematic low frequency variability biases in GCM simu-
985 lations, *Water Resources Research*, 48, W12 601, <https://doi.org/10.1029/2012WR012446>, 2012.



- Mehrotra, R. and Sharma, A.: Correcting for systematic biases in multiple raw GCM variables across a range of timescales, *Journal of Hydrology*, 520, 214–223, <https://doi.org/10.1016/j.jhydrol.2014.11.037>, 2015.
- Mehrotra, R. and Sharma, A.: A multivariate quantile-matching bias correction approach with auto- and cross-dependence across multiple time scales: Implications for downscaling, *Journal of Climate*, 29, 3519–3539, <https://doi.org/10.1175/jcli-d-15-0356.1>, 2016.
- 990 Meyer, J., Kohn, I., Stahl, K., Hakala, K., Seibert, J., and Cannon, A. J.: Effects of univariate and multivariate bias correction on hydrological impact projections in alpine catchments, *Hydrology and Earth System Sciences*, 23, 1339–1354, <https://doi.org/10.5194/hess-23-1339-2019>, 2019.
- Mezzadri, F.: How to generate random matrices from the classical compact groups, *Notices of the American Mathematical Society*, 54, 592–604, 2007.
- 995 Michelangeli, P.-A., Vrac, M., and Loukos, H.: Probabilistic downscaling approaches: Application to wind cumulative distribution functions, *Geophysical Research Letters*, 36, L11 708, <https://doi.org/10.1029/2009GL038401>, 2009.
- Milly, P. C. D., Betancourt, J., Falkenmark, M., Hirsch, R. M., Kundzewicz, Z. W., Lettenmaier, D. P., and Stouffer, R. J.: Stationarity is dead: Whither water management?, *Science*, 319, 573–574, <https://doi.org/10.1126/science.1151915>, 2008.
- Moore, R. J.: The PDM rainfall-runoff model, *Hydrology and Earth System Sciences*, 11, 483–499, [https://doi.org/10.5194/hess-11-483-](https://doi.org/10.5194/hess-11-483-2007)
1000 2007, 2007.
- Nahar, J., Johnson, F., and Sharma, A.: Assessing the extent of non-stationary biases in GCMs, *Journal of Hydrology*, 549, 148–162, <https://doi.org/10.1016/j.jhydrol.2017.03.045>, 2017.
- Nelsen, R. B.: *An Introduction to Copulas*, 2nd, New York: Springer Science Business Media, 2006.
- Nguyen, H., Mehrotra, R., and Sharma, A.: Correcting for systematic biases in GCM simulations in the frequency domain, *Journal of*
1005 *Hydrology*, 538, 117–126, <https://doi.org/10.1016/j.jhydrol.2016.04.018>, 2016.
- Nguyen, H., Mehrotra, R., and Sharma, A.: Correcting systematic biases across multiple atmospheric variables in the frequency domain, *Climate Dynamics*, 52, 1283–1298, <https://doi.org/10.1007/s00382-018-4191-6>, 2018.
- Olsson, J., Berggren, K., Olofsson, M., and Viklander, M.: Applying climate model precipitation scenarios for urban hydrological assessment: A case study in Kalmar City, Sweden, *Atmospheric Research*, 92, 364–375, <https://doi.org/10.1016/j.atmosres.2009.01.015>, 2009.
- 1010 Panofsky, H. A., Brier, G. W., and Best, W. H.: *Some Application of Statistics to Meteorology*, Earth and Mineral Sciences Continuing Education, College of Earth and Mineral Sciences, Pennsylvania State University, 1958.
- Papalexiou, S. M. and Montanari, A.: Global and regional increase of precipitation extremes under global warming, *Water Resources Research*, 55, 4901–4914, <https://doi.org/10.1029/2018WR024067>, 2019.
- Penman, H. L.: Natural evaporation from open water, bare soil and grass, *Proc. R. Soc. Lond. A*, 193, 120–145, 1948.
- 1015 Peyré, G. and Cuturi, M.: *Computational Optimal Transport*, vol. 11, Now Publishers, <https://doi.org/10.1561/22000000073>, 2019.
- Pham, M. T.: Copula-based stochastic modelling of evapotranspiration time series conditioned on rainfall as design tool in water resources management, PhD thesis, Faculty of Biosciences Engineering, Ghent University, 2016.
- Pham, M. T., Vernieuwe, H., De Baets, B., and Verhoest, N.: A coupled stochastic rainfall–evapotranspiration model for hydrological impact analysis, *Hydrology and Earth System Sciences*, 22, 1263–1283, <https://doi.org/10.5194/hess-22-1263-2018>, 2018.
- 1020 Piani, C. and Haerter, J. O.: Two dimensional bias correction of temperature and precipitation copulas in climate models, *Geophysical Research Letters*, 39, <https://doi.org/10.1029/2012gl053839>, 2012.
- Piani, C., Haerter, J. O., and Coppola, E.: Statistical bias correction for daily precipitation in regional climate models over Europe, *Theoretical and Applied Climatology*, 99, 187–192, <https://doi.org/10.1007/s00704-009-0134-9>, 2010.



- Pitié, F., Kokaram, A. C., and Dahyot, R.: Automated colour grading using colour distribution transfer, *Computer Vision and Image Understanding*, 107, 123–137, <https://doi.org/10.1016/j.cviu.2006.11.011>, 2007.
- 1025 Popke, D., Stevens, B., and Voigt, A.: Climate and climate change in a radiative-convective equilibrium version of ECHAM6, *Journal of Advances in Modeling Earth Systems*, 5, 1–14, <https://doi.org/10.1029/2012MS000191>, 2013.
- Prechelt, L.: Automatic early stopping using cross validation: quantifying the criteria, *Neural Networks*, 11, 761–767, [https://doi.org/10.1016/S0893-6080\(98\)00010-0](https://doi.org/10.1016/S0893-6080(98)00010-0), 1998.
- 1030 Prein, A. F., Langhans, W., Fosser, G., Ferrone, A., Ban, N., Goergen, K., Keller, M., Tölle, M., Gutjahr, O., Feser, F., Brisson, E., Kollet, S., Schmidli, J., Lipzig, N. P. M., and Leung, R.: A review on regional convection-permitting climate modeling: Demonstrations, prospects, and challenges, *Reviews of Geophysics*, 53, 323–361, <https://doi.org/10.1002/2014RG000475>, 2015.
- Rajczak, J., Kotlarski, S., and Schär, C.: Does quantile mapping of simulated precipitation correct for biases in transition probabilities and spell lengths?, *Journal of Climate*, 29, 1605–1615, <https://doi.org/10.1175/JCLI-D-15-0162.1>, 2016.
- 1035 Rätty, O., Räisänen, J., and Ylhäisi, J. S.: Evaluation of delta change and bias correction methods for future daily precipitation: intermodel cross-validation using ENSEMBLES simulations, *Climate Dynamics*, 42, 2287–2303, <https://doi.org/10.1007/s00382-014-2130-8>, 2014.
- Rätty, O., Räisänen, J., Bosshard, T., and Donnelly, C.: Intercomparison of univariate and joint bias correction methods in changing climate from a hydrological perspective, *Climate*, 6, 33, <https://doi.org/10.3390/cli6020033>, 2018.
- Reiter, P., Gutjahr, O., Schefczyk, L., Heinemann, G., and Casper, M.: Does applying quantile mapping to subsamples improve the bias correction of daily precipitation?, *International Journal of Climatology*, 38, 1623–1633, <https://doi.org/10.1002/joc.5283>, 2018.
- 1040 Rizzo, M. L. and Székely, G. J.: Energy distance, *Wiley Interdisciplinary Reviews: Computational Statistics*, 8, 27–38, <https://doi.org/10.1002/wics.1375>, 2016.
- Robin, Y., Vrac, M., Naveau, P., and Yiou, P.: Multivariate stochastic bias corrections with optimal transport, *Hydrology and Earth System Sciences*, 23, 773–786, <https://doi.org/10.5194/hess-23-773-2019>, 2019.
- 1045 Rojas, R., Feyen, L., Dosio, A., and Bavera, D.: Improving pan-European hydrological simulation of extreme events through statistical bias correction of RCM-driven climate simulations, *Hydrology & Earth System Sciences*, 15, <https://doi.org/10.5194/hess-15-2599-2011>, 2011.
- Salas, J. D.: *Applied Modeling of Hydrologic Time Series*, Water Resources Publication, 1980.
- Schmidli, J., Frei, C., and Vidale, P. L.: Downscaling from GCM precipitation: a benchmark for dynamical and statistical downscaling methods, *International Journal of Climatology*, 26, 679–689, <https://doi.org/10.1002/joc.1287>, 2006.
- 1050 Schölzel, C. and Friederichs, P.: Multivariate non-normally distributed random variables in climate research-introduction to the copula approach, *Nonlinear Processes in Geophysics*, 15, 761–772, <https://doi.org/10.5194/npg-15-761-2008>, 2008.
- Shepherd, T. G.: Atmospheric circulation as a source of uncertainty in climate change projections, *Nature Geoscience*, 7, 703–708, <https://doi.org/10.1038/ngeo2253>, 2014.
- 1055 Sippel, S., Otto, F. E. L., Forkel, M., Allen, M. R., Guillod, B. P., Heimann, M., Reichstein, M., Seneviratne, S. I., Thonicke, K., and Mahecha, M. D.: A novel bias correction methodology for climate impact simulations, *Earth System Dynamics*, 7, 71–88, <https://doi.org/10.5194/esd-7-71-2016>, 2016.
- Srikanthan, R. and Pegram, G. G. S.: A nested multisite daily rainfall stochastic generation model, *Journal of Hydrology*, 371, 142–153, <https://doi.org/10.1016/j.jhydrol.2009.03.025>, 2009.
- 1060 Strandberg, G., Barring, L., Hansson, U., Jansson, C., Jones, C., Kjellström, E., Kupiainen, M., Nikulin, G., Samuelsson, P., and Ullerstig, A.: CORDEX scenarios for Europe from the Rossby Centre regional climate model RCA4, Tech. rep., SMHI, 2015.



- Sunyer, M. A., Madsen, H., Rosbjerg, D., and Arnbjerg-Nielsen, K.: A Bayesian approach for uncertainty quantification of extreme precipitation projections including climate model interdependency and nonstationary bias, *Journal of Climate*, 27, 7113–7132, <https://doi.org/10.1175/JCLI-D-13-00589.1>, 2014.
- 1065 Székely, G. J. and Rizzo, M. L.: Testing for equal distributions in high dimension, *InterStat*, 5, 1249–1272, 2004.
- Székely, G. J. and Rizzo, M. L.: Energy statistics: A class of statistics based on distances, *Journal of statistical planning and inference*, 143, 1249–1272, <https://doi.org/10.1016/j.jspi.2013.03.018>, 2013.
- Teutschbein, C. and Seibert, J.: Bias correction of regional climate model simulations for hydrological climate-change impact studies: Review and evaluation of different methods, *Journal of Hydrology*, 456, 12–29, <https://doi.org/10.1016/j.jhydrol.2012.05.052>, 2012.
- 1070 Teutschbein, C. and Seibert, J.: Is bias correction of regional climate model (RCM) simulations possible for non-stationary conditions?, *Hydrology and Earth System Sciences*, 17, 5061–5077, <https://doi.org/10.5194/hess-17-5061-2013>, 2013, 2013.
- Themeßl, M. J., Gobiet, A., and Leuprecht, A.: Empirical-statistical downscaling and error correction of daily precipitation from regional climate models, *International Journal of Climatology*, 31, 1530–1544, <https://doi.org/10.1002/joc.2168>, 2011.
- Themeßl, M. J., Gobiet, A., and Heinrich, G.: Empirical-statistical downscaling and error correction of regional climate models and its impact on the climate change signal, *Climatic Change*, 112, 449–468, <https://doi.org/10.1007/s10584-011-0224-4>, 2012.
- 1075 Tokarska, K. B., Stolpe, M. B., Sippel, S., Fischer, E. M., Smith, C. J., Lehner, F., and Knutti, R.: Past warming trend constrains future warming in CMIP6 models, *Science Advances*, 6, eaaz9549, <https://doi.org/10.1126/sciadv.aaz9549>, 2020.
- Van Schaeybroeck, B. and Vannitsem, S.: Assessment of calibration assumptions under strong climate changes, *Geophysical Research Letters*, 43, 1314–1322, <https://doi.org/10.1002/2016GL067721>, 2016.
- 1080 Van Vuuren, D. P., Edmonds, J., Kainuma, M., Riahi, K., Thomson, A., Hibbard, K., Hurtt, G. C., Kram, T., Krey, V., Lamarque, J.-F., et al.: The representative concentration pathways: an overview, *Climatic Change*, 109, 5, <https://doi.org/10.1007/s10584-011-0148-z>, 2011.
- Vandenbergh, S., Verhoest, N. E. C., Onof, C., and De Baets, B.: A comparative copula-based bivariate frequency analysis of observed and simulated storm events: A case study on Bartlett-Lewis modeled rainfall, *Water Resources Research*, 47, W07 529, 2011.
- Velázquez, J. A., Troin, M., Caya, D., and Brissette, F.: Evaluating the time-invariance hypothesis of climate model bias correction: implications for hydrological impact studies, *Journal of Hydrometeorology*, 16, 2013–2026, <https://doi.org/10.1175/JHM-D-14-0159.1>, 2015.
- 1085 Verhoest, N., Troch, P. A., and De Troch, F. P.: On the applicability of Bartlett–Lewis rectangular pulses models in the modeling of design storms at a point, *Journal of Hydrology*, 202, 108–120, [https://doi.org/10.1016/S0022-1694\(97\)00060-7](https://doi.org/10.1016/S0022-1694(97)00060-7), 1997.
- Verstraeten, G., Poesen, J., Demarée, G., and Salles, C.: Long-term (105 years) variability in rain erosivity as derived from 10-min rainfall depth data for Ukkel (Brussels, Belgium): Implications for assessing soil erosion rates, *Journal of Geophysical Research*, 111, D22 109, <https://doi.org/10.1029/2006jd007169>, 2006.
- 1090 Villani, C.: *Optimal transport: old and new*, vol. 338, Springer Science & Business Media, 2008.
- Vrac, M.: Multivariate bias adjustment of high-dimensional climate simulations: the Rank Resampling for Distributions and Dependences (R2D2) bias correction, *Hydrology and Earth System Sciences*, 22, 3175, <https://doi.org/10.5194/hess-22-3175-2018>, 2018.
- Vrac, M. and Friederichs, P.: Multivariate—intervariable, spatial, and temporal—bias correction, *Journal of Climate*, 28, 218–237, <https://doi.org/10.1175/JCLI-D-14-00059.1>, 2015.
- 1095 Wang, L. and Chen, W.: Equiratio cumulative distribution function matching as an improvement to the equidistant approach in bias correction of precipitation, *Atmospheric Science Letters*, 15, 1–6, <https://doi.org/10.1002/asl2.454>, 2014.
- Wang, Y., Sivandran, G., and Bielicki, J. M.: The stationarity of two statistical downscaling methods for precipitation under different choices of cross-validation periods, *International Journal of Climatology*, 38, e330–e348, <https://doi.org/10.1002/joc.5375>, 2018.



- 1100 Wilcke, R. A. I., Mendlik, T., and Gobiet, A.: Multi-variable error correction of regional climate models, *Climatic Change*, 120, 871–887, <https://doi.org/10.1007/s10584-013-0845-x>, 2013.
- Willems, P.: Revision of urban drainage design rules after assessment of climate change impacts on precipitation extremes at Uccle, Belgium, *Journal of Hydrology*, 496, 166–177, <https://doi.org/10.1016/j.jhydrol.2013.05.037>, 2013.
- Willems, P. and Vrac, M.: Statistical precipitation downscaling for small-scale hydrological impact investigations of climate change, *Journal*
- 1105 of *Hydrology*, 402, 193–205, <https://doi.org/10.1016/j.jhydrol.2011.02.030>, 2011.
- Yang, W., Gardelin, M., Olsson, J., and Bosshard, T.: Multi-variable bias correction: application of forest fire risk in present and future climate in Sweden, *Natural Hazards and Earth System Sciences*, 15, 2037–2057, <https://doi.org/10.5194/nhess-15-2037-2015>, 2015.
- Zhang, X., Alexander, L., Hegerl, G. C., Jones, P., Tank, A. K., Peterson, T. C., Trewin, B., and Zwiers, F. W.: Indices for monitoring changes in extremes based on daily temperature and precipitation data, *Wiley Interdisciplinary Reviews: Climate Change*, 2, 851–870,
- 1110 <https://doi.org/10.1002/wcc.147>, 2011.
- Zscheischler, J., Westra, S., Hurk, B. J. J. M., Seneviratne, S. I., Ward, P. J., Pitman, A., AghaKouchak, A., Bresch, D. N., Leonard, M., Wahl, T., and Zhang, X.: Future climate risk from compound events, *Nature Climate Change*, p. 1, <https://doi.org/10.1038/s41558-018-0156-3>, 2018.
- Zscheischler, J., Fischer, E. M., and Lange, S.: The effect of univariate bias adjustment on multivariate hazard estimates, *Earth System*
- 1115 *Dynamics*, 10, 31–43, <https://doi.org/10.5194/esd-10-31-2019>, 2019.
- Zscheischler, J., Martius, O., Westra, S., Bevacqua, E., Raymond, C., Horton, R. M., van den Hurk, B., AghaKouchak, A., Jézéquel, A., Mahecha, M. D., Maraun, D., Ramos, A. M., Ridder, N. N., Thiery, W., and Vignotto, E.: A typology of compound weather and climate events, *Nature Reviews Earth & Environment*, <https://doi.org/10.1038/s43017-020-0060-z>, 2020.