

Appendix A: Analytical approximation of expectation value for the normalisation-induced bias

Assumptions and Notation:

- Assume independent and identically distributed (i.e., stationary) variables $X_{t,i}$ with mean given by $\mathbf{E}(X) = \mu$ and variance $\mathbf{Var}(X) = \sigma^2$. Let the subscripts t and i denote time and grid cell index, respectively.
- 5 – Denote t as an arbitrary time step in the ‘out-of-base’ (independent) period, and t_{ref} as an arbitrary time step inside the reference period. Let n_{ref} denote the length of the reference period.
- Denote $\Delta_{bias} = \mathbf{E}\left(\frac{X_{t,i}}{\hat{\mu}_{ref,i}}\right) - 1$ as the change induced by normalisation by the mean of an independent reference period (i.e., ‘normalisation bias’).

Our objective is to find an analytical approximation of the expectation value for the artificially induced relative change
 10 (‘bias’) by dividing a random variable $X_{t,i}$ as defined above by its sample mean ($\hat{\mu}_{ref,i} = \frac{1}{n} \sum_{t_{ref}=1}^{n_{ref}} X_{t_{ref},i}$, where $\mathbf{E}(\hat{\mu}_{ref,i}) = \mu$) that has been estimated in an independent (‘reference’) period, i.e.

$$\Delta_{bias} = \mathbf{E}\left(\frac{X_{t,i}}{\hat{\mu}_{ref,i}}\right) - 1 \approx f(\mu, \sigma, n_{ref}). \quad (\text{A1})$$

Clearly, an unbiased estimate of this quantity in stationary time series should yield $\Delta_{bias} = 0$. Because $X_{t,i}$ and $\hat{\mu}_{ref,i}$ are independent, we can write,

$$15 \quad \Delta_{bias} = \mathbf{E}(X_{t,i})\mathbf{E}\left(\frac{1}{\hat{\mu}_{ref,i}}\right) - 1 = \mu\mathbf{E}\left(\frac{1}{\hat{\mu}_{ref,i}}\right) - 1. \quad (\text{A2})$$

If we substitute $\hat{\mu}_{ref,i} = \mu(1 + \epsilon_{ref,i})$, where $\mathbf{E}(\epsilon_i) = 0$, $\mathbf{Var}(\epsilon_i) = \frac{\sigma^2}{\mu^2 n_{ref}}$, and the subscript ref has been dropped from ϵ_i for convenience, we get

$$\Delta_{bias} = \mu\mathbf{E}\left(\frac{1}{\mu(1 + \epsilon_i)}\right) - 1 = \mathbf{E}\left(\frac{1}{1 + \epsilon_i}\right) - 1. \quad (\text{A3})$$

A Taylor expansion around the function $g(x) = \frac{1}{1+x}$ at $x = 0$ yields

$$20 \quad g(x) = \frac{1}{1+x} = 1 - x + x^2 - x^3 + x^4 - x^5 + \dots \quad (\text{A4})$$

We will see below that the convergence criterion $\epsilon_i < |1|$ of the Taylor series is met in practically relevant cases, but it should be noted that convergence cannot be ensured in all theoretically conceivable cases. Using Taylor expansion, Δ_{bias} can be approximated, making use of the linearity of the expectation operator $\mathbf{E}()$ and of the fact that $\mathbf{E}(\epsilon_i) = 0$ and $\mathbf{E}(\epsilon_i^2) = \mathbf{Var}(\epsilon_i^2) = \frac{\sigma^2}{\mu^2 n_{ref}}$ by definition,

$$25 \quad \Delta_{bias} = \mathbf{E}\left(\frac{1}{1 + \epsilon_i}\right) - 1 \quad (\text{A5})$$

$$= \mathbf{E}(1 - \epsilon_i + \epsilon_i^2 - \epsilon_i^3 + \epsilon_i^4 - \epsilon_i^5 + \dots) - 1 \quad (\text{A6})$$

$$= \frac{\sigma^2}{\mu^2 n_{ref}} - \mathbf{E}(\epsilon_i^3) + \mathbf{E}(\epsilon_i^4) - \mathbf{E}(\epsilon_i^5) + \dots \quad (\text{A7})$$

This expression yields a sum over the central moments of the distribution of ϵ_i 's. For a symmetric probability distribution (recall that ϵ_i denote the deviations of the sample means in a reference period around the underlying true mean), because $E(\epsilon_i^k) = 0$, where k is any uneven exponent $k \in \mathbb{N}$, Eq. A7 reduces further to

$$\Delta_{bias} = \frac{\sigma^2}{\mu^2 n_{ref}} + \mathbf{E}(\epsilon_i^4) + \mathbf{E}(\epsilon_i^6) + \dots \quad (\text{A8})$$

- 5 As long as $\epsilon_i < |1|$ is fulfilled, the quadratical term dominates both Eq. A7 and Eq. A8. The present analytical approximation (both Eq. A7 and Eq. A8) provides the important insights that 1) normalisation with a 'reference period mean' leads to an artificial increase of spatial averages in the out-of-base period, i.e. the bias is always positive in the out-of-base period, $\Delta_{bias} > 0$, and 2) that $\Delta_{bias} \propto (\frac{\sigma}{\mu} \frac{1}{\sqrt{n_{ref}}})^2$, i.e. the coefficient of variation in the distribution of sample means (i.e., $c_v[\hat{\mu}_{ref,i}] = \frac{\sigma}{\mu\sqrt{n_{ref}}}$). For any fixed n_{ref} , the amplitude of the normalisation-induced biases only depends on the ratio $\frac{\sigma}{\mu}$. We verify below numerically
- 10 that this approximation works well for random variables $X_{t,i}$ drawn from i. a Gaussian distribution, ii. a Generalized Extreme Value distribution with two different choices for the shape parameter ($\xi = 0$, 'Gumbel distribution', and $\xi \neq 0$).

Gaussian distribution

Assume $X_{t,i} \sim \mathcal{N}(\mu, \sigma^2)$, the sample mean deviations from the true mean will follow $\epsilon_i \sim \mathcal{N}(0, \frac{\sigma^2}{\mu^2 n_{ref}})$. If we substitute with $\epsilon_i = \frac{\sigma}{\mu} \frac{1}{\sqrt{n_{ref}}} Y$, where $Y \sim \mathcal{N}(0, 1)$ in Eq. A8, the above expression reduces to

$$15 \quad \Delta_{bias} = \frac{\sigma^2}{\mu^2 n_{ref}} + (\frac{\sigma}{\mu} \frac{1}{\sqrt{n_{ref}}})^4 \mathbf{E}(Y^4) + (\frac{\sigma}{\mu} \frac{1}{\sqrt{n_{ref}}})^6 \mathbf{E}(Y^6) + \dots \quad (\text{A9})$$

Because higher-order moments of a standard normal distributed random variable are well-known and can be derived analytically (Johnson et al., 1994, i.e., $\mathbf{E}(Y^4) = 3$, $\mathbf{E}(Y^6) = 15$), an analytical expression of the normalisation-induced bias becomes straightforward:

$$\Delta_{bias} \approx \frac{\sigma^2}{\mu^2 n_{ref}} + 3(\frac{\sigma}{\mu} \frac{1}{\sqrt{n_{ref}}})^4 + 15(\frac{\sigma}{\mu} \frac{1}{\sqrt{n_{ref}}})^6. \quad (\text{A10})$$

- 20 A comparison of Eq. A10 (i.e. the first three terms in the Taylor approximation) to numerical simulations shows that the analytical approximation works well (Fig. A-1a). Furthermore, the estimation of mean and standard deviation from the empirical time series to calculate the expectation value for the biases is unbiased and show surprisingly little variation (Fig. A-1b) even for a relatively small number of grid cells, where random variation in stationary time series becomes considerable (Fig. A-1b).

- However, one important caveat is that Eq. A3 and the subsequent approximation only works as long as $\epsilon_i < |1|$ is fulfilled.
- 25 How likely is a violation of this criterion? Numerical simulations for $n_{ref} = 30$ appear to be very stable for any $\frac{\sigma}{\mu} > 0.8$ in the $X_{t,i}$'s, i.e. corresponding roughly to a $C_v[\hat{\mu}_{ref,i}] \approx 0.2$. This would mean that the chance of $|\epsilon_i| \geq 1$ corresponds to a -5σ event with a probability of roughly 1 to 3.5 million. Given that the observed $\frac{\sigma}{\mu}$ ratios are considerably larger than the values tested here even in the driest regions of the world, we conclude that the approximation can be used for the vast majority, if not all, practical considerations.

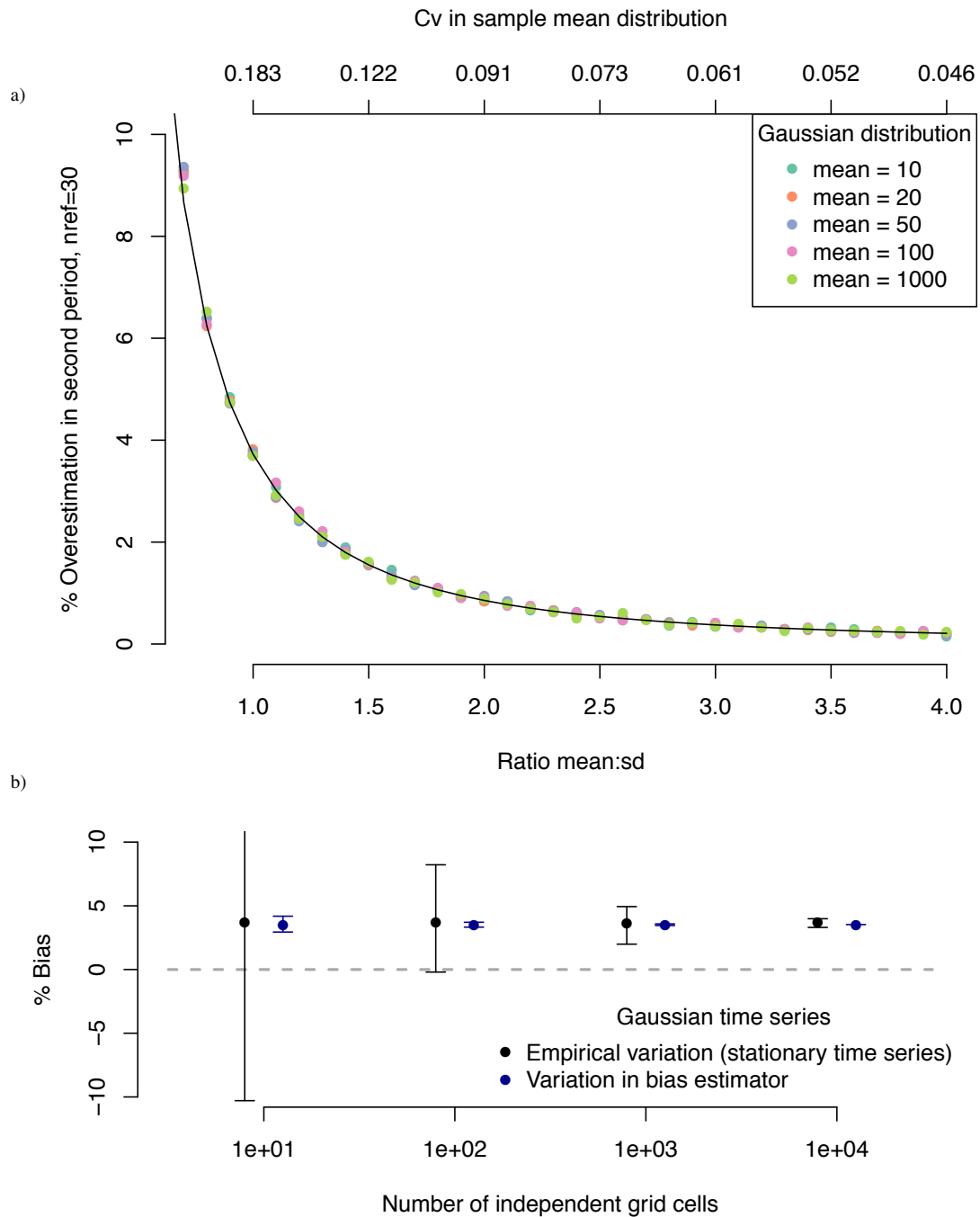


Figure A-1. a. Ratio of mean to sd vs. normalisation-induced bias in a Gaussian distribution for numerical simulations with various mean values (dots), and the derived analytical approximation. The reference period length is taken as $n_{ref} = 30$, and numerical simulations are conducted with $n = 10^5$ grid cells with each 60 time steps. b. Variation in the empirical estimates of the biases (darkblue) for a given number of independent grid cells ($\frac{\mu}{\sigma} = 1$, $n_{ref} = 30$). The magnitude of random changes in stationary time series with $n_{ref} = 30$ and $n_{obase} = 30$ is shown for comparison in black. Error bars indicate 5th and 95th percentile in repeated numerical simulations.

Generalized extreme value distribution

An important question that arises from the expression outlined in Eq. A7 is whether the higher-order terms in the Taylor approximation can be ignored in practical applications, where an assumption of Gaussianity might not hold. Here, we test this proposition for the Generalized Extreme Value distribution as an appropriate model for annual maxima as investigated in the main manuscript with two different choice for the distribution's shape parameter (ξ).

i. Gumbel distribution

Therefore, we first assume, in analogy to the paragraph above, independent and identically distributed (i.e., stationary) random variables drawn from a Generalized Extreme Value distribution with zero shape parameter ('Gumbel distribution', $X_{t,i} \sim GEV(\mu', \sigma', \xi = 0)$, where μ' , σ' and $\xi = 0$ denote the GEV's location, scale and shape parameter, respectively Johnson et al., 1995). The expectations for mean (μ) and variance (σ^2) of a GEV are given by $\mu = \mu' + \sigma'\gamma$, where γ denotes Euler's constant.

Following Eq. A7, we can readily derive an analytical expression for the expectation value of the normalisation-induced bias, i.e.

$$\Delta_{bias} = \frac{\sigma^2}{\mu^2 n_{ref}} - \mathbf{E}(\epsilon_i^3) + \mathbf{E}(\epsilon_i^4) - \mathbf{E}(\epsilon_i^5) + \dots \quad (\text{A11})$$

$$= \left(\frac{\pi}{\sqrt{6n_{ref}} \left(\frac{\mu'}{\sigma'} + \gamma \right)} \right)^2 - \mathbf{E}(\epsilon_i^3) + \mathbf{E}(\epsilon_i^4) - \mathbf{E}(\epsilon_i^5) + \dots \quad (\text{A12})$$

Here, we note again that the quadratical term dominates the expression. If we make the simplifying assumption that the sample means $\hat{\mu}_{ref,i}$ for $n_{ref} = 30$ follow (approximately) a Gaussian distribution (the assumption is only needed for the higher order terms of the Taylor expansion), we can derive an analytical approximation for the normalisation-induced bias by insertion and in analogy to above, i.e.

$$\Delta_{bias} \approx \left(\frac{\pi}{\sqrt{6n_{ref}} \left(\frac{\mu'}{\sigma'} + \gamma \right)} \right)^2 + \left(\frac{\sigma}{\mu} \frac{1}{\sqrt{n_{ref}}} \right)^4 \mathbf{E}(Y^4) + \dots \quad (\text{A13})$$

$$\approx \left(\frac{\pi}{\sqrt{6n_{ref}} \left(\frac{\mu'}{\sigma'} + \gamma \right)} \right)^2 + 3 \left(\frac{\pi}{\sqrt{6n_{ref}} \left(\frac{\mu'}{\sigma'} + \gamma \right)} \right)^4. \quad (\text{A14})$$

Similarly to above, we find that the ratio of location to scale parameter ($\frac{\mu'}{\sigma'}$), for any fixed reference period length (n_{ref}), determines the magnitude of the bias. The analytical approximation can be verified by numerical simulation using GEV-distributed random variables, and is found to fit the data very well (Fig. A-2a). Furthermore, an estimator of the expectation value of the biases by only estimating the mean and standard deviation of empirical time series (i.e., using the first term in the Taylor approximation) can be derived easily and is found to work reliable also for a small number of independent grid cells (Fig. A-2c).

ii. GEV distribution with $\xi \neq 0$

Here, we test whether the analytical argument from above can be extended to Generalized Extreme Value distributions with $\xi \neq 0$. In practical applications of the GEV to observed maximum precipitation, a shape parameter of $\xi \approx 0.2$ is often found

(Van den Brink and Können, 2011), therefore we test here for $X_{t,i} \sim \text{GEV}(\mu', \sigma', \xi = 0.2)$. The expectations for mean (μ) and variance (σ^2) of a GEV, when $0 > \epsilon < 1$, are given by $\mu = \mu' + \sigma' \frac{\Gamma(1-\xi)-1}{\xi}$ and $\sigma^2 = (\sigma')^2 \frac{(g_2 - g_1^2)}{\xi}$, where $g_k = \Gamma(1 - k\xi)$, $k = 1, 2$, and $\Gamma(t)$ is the gamma function (Johnson et al., 1995).

Hence, the (dominant) quadratic term in the Taylor approximation in Eq. A7 reads,

$$5 \quad \Delta_{bias} \approx \frac{(g_2 - g_1^2)}{n_{\text{ref}} \xi \left[\frac{\mu'}{\sigma'} + \frac{\Gamma(1-\xi)-1}{\xi} \right]^2}. \quad (\text{A15})$$

The approximation works again very well in numerical simulations (Fig. A-2b), and shows that if $\xi \neq 0$, there is a dependency on ξ , n_{ref} , and again the ratio of $\frac{\mu'}{\sigma'}$ (rather than either μ' or σ' individually), which determine the magnitude of the normalisation-induced bias. Please note that the approximation works similarly well for random variables drawn from a GEV-distribution with negative shape parameter ($\xi = -0.2$, not shown).

10 Short Remark on in-stationarities in the out-of-base period

Many real-world precipitation time series show in-stationarities due to climatic variations (O’Gorman, 2015) that are typically diagnosed as relative changes in the precipitation amount. Hence, the question whether and how any ‘real change in the expectation value’ outside the reference period can be disentangled from normalisation-induced biases becomes topical. Given the analytical approximation above, it becomes obvious that the highlighted normalisation-induced bias scales in-stationarities

15 in the out-of-base period in a multiplicative way:

Let c denote any change between the reference period expectation and some future period (e.g. assume one is interested in global or latitudinal changes in a past and future climatic period), i.e. such that $\mathbf{E}(X_{t_{\text{ref}},i}) = c\mathbf{E}(X_{t,i})$, then the bias (Δ_{bias} , after accounting for the ‘real change’) would simply scale with the relative change (Δ denotes the total apparent change):

$$\Delta = c\mathbf{E}\left(\frac{X_{t,i}}{\hat{\mu}_{\text{ref},i}}\right) - 1 \quad (\text{A16})$$

$$20 \quad = c\mathbf{E}\left(\frac{1}{1 + \epsilon_i}\right) - 1 \quad (\text{A17})$$

$$= \underbrace{c - 1}_{\text{true change}} + c \underbrace{\left[\frac{\sigma^2}{\mu^2 n_{\text{ref}}} - \mathbf{E}(\epsilon_i^3) + \mathbf{E}(\epsilon_i^4) - \mathbf{E}(\epsilon_i^5) + \dots \right]}_{\Delta_{\text{bias}}} \quad (\text{A18})$$

From Eq. A18, it is straightforward to see that for any multiplicative changes in the expectation of the out-of-base variables, the normalisation-induced bias scales with the change. Hence, this expression implies that to detect the ‘true change c ’ between two periods, the normalisation-induced bias has to be accounted for, i.e.

$$25 \quad c = \frac{\Delta + 1}{1 + \Delta_{\text{bias}}}. \quad (\text{A19})$$

Numerical simulations can be easily conducted similar to Subsection 1.1 and 1.2 to verify that this scaling holds (not shown).

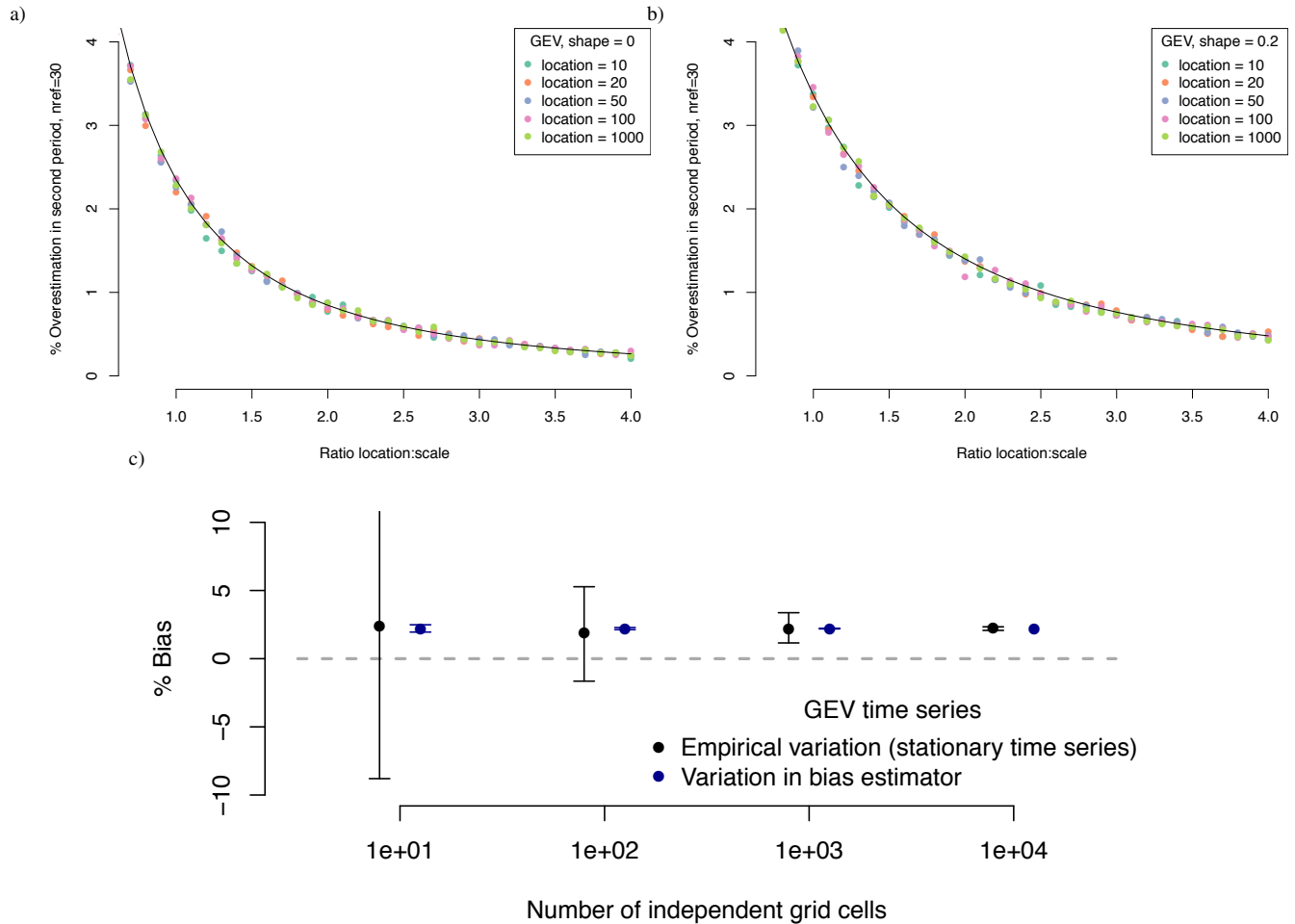


Figure A-2. a) Ratio of location to scale parameter vs. normalisation-induced bias in a Generalized extreme value distribution for numerical simulations with various location parameter values (dots) and a) zero shape parameter, and b) with $\xi = 0.2$. Reference period length is taken as $n_{ref} = 30$, and numerical simulations are conducted with $n = 10^5$ grid cells with each 60 time steps. c) Variation in the empirical estimates of the biases (darkblue) for a given number of independent grid cells ($\frac{\mu'}{\sigma'} = 1$, $\xi = 0$, $n_{ref} = 30$). The magnitude of random changes in stationary time series with $n_{ref} = 30$ and $n_{obase} = 30$ is shown for comparison in black. Error bars indicate 5th and 95th percentile in repeated numerical simulations.