

1           **Sensitivity of the SIMulation-EXtrapolation (SIMEX)**  
2           **Methodology to Mis-specification of the Statistical Properties**  
3                       **of the Measurement Errors**

4  
5           GABRIELE VILLARINI<sup>1</sup>, DARIO TREPPIEDI<sup>2</sup>, LEONARDO V. NOTO<sup>2</sup>

6  
7  
8           <sup>1</sup>IIHR—Hydroscience and Engineering, University of Iowa, Iowa City, USA

9           <sup>2</sup>Dipartimento di Ingegneria, Università degli Studi di Palermo, Palermo, Italy

10  
11  
12                                       Manuscript submitted to

13                                       *Theoretical and Applied Climatology*

14                                       2 September 2022

15  
16                                       *Revised February 2023*

17  
18  
19           *Corresponding author address:*

20           Gabriele Villarini, IIHR-Hydroscience & Engineering, The University of Iowa, 323B C.

21           Maxwell Stanley Hydraulics Laboratory, Iowa City, 52242, Iowa, USA. E-mail: gabriele-

22           villarini@uiowa.edu. Tel.: (319) 384-0596

23  
24

1 ABSTRACT

2  
3 In hydrometeorological and environmental studies, it is common to seek relations  
4 between two variables (predictand and predictor), one of which (predictor) is affected by  
5 uncertainties. These errors unavoidably affect the results of the analyses by providing  
6 erroneous estimates of the parameters of the predictor-predictand model. A possible  
7 solution is represented by the SIMulation-EXtrapolation (SIMEX) methodology. This  
8 approach follows two steps: 1) perturbation of the predictor with increasing level of  
9 uncertainties (multiples of the known error variance); and 2) finding a relation between the  
10 model's parameters and level of uncertainty, which allows their extrapolation to the error-  
11 free case.

12 The application of the SIMEX methodology requires the a priori knowledge of the  
13 mean, variance, and distribution of the measurement errors. However, in hydrologic and  
14 climatologic studies, this is not the case and the impact of an erroneous specification of  
15 these statistical properties on the results of the analyses has received little attention. The  
16 aim of this study is to investigate the sensitivity of the SIMEX methodology to mis-  
17 specification of the error characteristics. By using a simulation-based approach, we  
18 investigate the impact of an imperfect knowledge of the characteristics of the errors  
19 associated with the predictor (mean, variance, and probability distribution). Our results  
20 suggest that SIMEX is robust against mis-specification of the moments and distribution of  
21 the measurement errors, that it performs better than standard linear regression, even when  
22 these statistical properties are erroneously specified and that, for these reasons, it could  
23 find an useful application to seasonal forecasting of hydroclimatic variables.

24  
25 Keywords: Simulation-extrapolation; SIMEX; Uncertainties; Measurement Errors

# 1 **1. Introduction**

2 In several disciplines, we are interested in finding the relation between two random  
3 variables, denoted by  $X$  and  $Y$ , where the former is the predictor and the latter the  
4 predictand. In general, while we usually assume that  $Y$  is error-free, this may not be the  
5 case for  $X$ : instead of measuring  $X$  (latent variable), we actually measure another variable,  
6  $W$ , which represents the combination of  $X$  and the measurement errors associated with it.  
7 In hydrometeorological and environmental studies, this is a very common problem. For  
8 instance, consider the case in which we want to relate radar rainfall estimates to the average  
9 of different rain gage measurements within the pixel of interest (considered to be the  
10 “ground truth”). When we relate these two quantities, we generally neglect to account for  
11 the uncertainties associated with computing pixel-averaged rainfall based on a limited  
12 number of rain gages (e.g., Villarini et al. 2008). As mentioned in Hasan et al. (2014), by  
13 neglecting these uncertainties, we introduce a bias when converting radar reflectivity to  
14 rainfall. Another example is related to the impacts of uncertainties in input variables in the  
15 estimation of the parameters of hydrologic models (e.g., Chowdhury and Sharma 2008).  
16 Moreover, consider for instance a time series of precipitation from a global climate model  
17 (GCM) that we want to fit with a gamma distribution. It is well-known that there are large  
18 uncertainties associated with GCM outputs, especially for a variable like precipitation. If  
19 we neglect to account for these uncertainties and apply standard approaches, our estimation  
20 of the parameters of the gamma distribution will be biased (Woldemeskel et al. 2014).

21 More specifically, let us focus on the case in which  $X$  and  $Y$  are linearly related:

$$22 \quad \mathbf{Y} = \beta_0 + \beta_1 \mathbf{X} + \varepsilon \quad (1)$$

1 where  $\varepsilon$  represents the uncertainties with respect to the regression line and is normally  
2 distributed with mean equal to zero and standard deviation equal to  $\sigma_\varepsilon$ .

3 Assuming a classic error model (Carroll et al. 2007), instead of measuring X we  
4 measure W, which accounts for the measurement errors U in X in an additive form:

$$5 \quad W=X+U \tag{2}$$

6 where U is independent of both X and Y and is described by a Gaussian distribution with  
7 mean  $\mu_U$  equal to zero and standard deviation  $\sigma_U$ .

8 Considering Y as the predictand and W as the predictor, ordinary least squares method  
9 (OLS) provides a consistent estimation of  $\hat{\beta}_{1,naive} = \lambda\beta_1$  rather than  $\beta_1$  (e.g., Carroll et al.  
10 2007) where:

$$11 \quad \lambda = \frac{\sigma_x^2}{\sigma_x^2 + \sigma_U^2} < 1 \tag{3}$$

12 and  $\sigma_x^2$  represents the variance of the variable X, and  $\lambda$  is called the reliability ratio (Fuller  
13 1987). The estimate of the slope by means of OLS is therefore biased when  $\sigma_U^2 > 0$ . In  
14 general, the presence of measurement errors not only introduces an attenuation of the slope  
15  $\beta_1$ , but makes the data noisier as well (e.g., Carroll et al. 2007).

16 Since it is clear how measurement errors may significantly affect the results of our  
17 analyses, they should be accounted for. In the field of applied statistics, this is a well-known  
18 problem and it has been the object of a vast literature (e.g., Fuller 1987; Gleser 1990;  
19 Brown and Mariano 1993; Carroll et al. 2007). Among all of the proposed approaches, a  
20 very effective method to reduce the bias associated with the presence of measurement  
21 errors is the SIMulation-EXtrapolation (SIMEX) method, first introduced by Cook and  
22 Stefanski (1994). SIMEX is based on the idea of adding measurement errors to the data as  
23 multiples of the known error variance and finding a relation between the targeted parameter

1 and the level of added noise, allowing the extrapolation of the results to the error-free case.  
2 We will discuss this approach in more details in the next section. Among other fields, this  
3 simulation-based method has been widely used in ecology (e.g., Kangas 1998 , Stoklosa et  
4 al. 2015; Ponzi et al. 2019; Kinane et al. 2021), and biostatistics and medicine (e.g., Lauzon  
5 et al. 2013; Guolo 2014; Alexeeff et al. 2016; Oh et al. 2018). However, in  
6 hydrometeorological and environmental studies, it has found limited application. The first  
7 application of this methodology was discussed by Chowdhury and Sharma (2007) to infill  
8 values of the Southern Oscillation Index based on sea surface temperature anomalies.  
9 Chowdhury and Sharma (2008) used the SIMEX methodology to quantify the bias in key  
10 storage parameters in the Sacramento Model. Then Woldemeskel et al (2012) showed the  
11 improvement in the estimation of the parameters of future droughts when SIMEX was used  
12 to account for uncertainties in the GCM outputs. Finally, Hasan et al. (2014) used this  
13 method when relating radar reflectivity to the ground truth, accounting for the uncertainties  
14 associated with a “ground truth” obtained from the average of multiple point  
15 measurements.

16 One of the possible obstacles towards a more widespread use of this methodology is  
17 the required knowledge of the statistical characteristics (i.e., mean, variance, and  
18 probability distribution) of the measurement error  $U$ . For the case of unknown variance,  
19 Devanarayan and Stefanski (2002) proposed a modification of the SIMEX approach, which  
20 however requires independent replicate measurements. However, as described in  
21 Chowdhury and Sharma (2008), an outstanding “issue to be investigated in greater detail  
22 is the specification of the error distribution for various hydro-climatological variables.”  
23 Therefore, what would happen if we mis-specify the error variance, or the mean, or even

1 the entire probability distribution still remains open questions. Carroll et al. (2007) wrote  
2 that “minor violations to the assumption of normality of the measurement errors is not  
3 critical in practice.” To the best of our knowledge, this is one of the very few indications  
4 about the robustness of SIMEX to mis-specification of the statistical characteristics of U.  
5 Therefore, a study in which this issue is addressed is still lacking and it would be important  
6 to show that even an imperfect knowledge of the measurement error characteristics could  
7 lead to more accurate results compared to the case in which we neglect them. In this study  
8 we tackle this question in a simulation framework, where we know the underlying relation  
9 between X and Y (Section 2). In Section 3 we will show that SIMEX tends to perform  
10 better than OLS even under mis-specification of the error characteristics, while Section 4  
11 discusses the main points made and closes the article.

## 12 **2. SIMulation-EXtrapolation: SIMEX**

13 There are several papers describing the SIMEX methodology in details for linear and  
14 non-linear models, homoschedastic or heteroschedastic errors, both in additive or  
15 multiplicative forms. The reader is pointed to Carroll et al. (2007) and references therein  
16 for more information. Here we consider the linear model in equation (1), and the additive  
17 error model in equation (2) to describe the measurement errors in X. We also assume that  
18 we know that U is normally distributed with mean equal to zero and variance equal to  $\sigma_U^2$ .

19 The SIMEX methodology consists of a simulation and an extrapolation step. In the  
20 simulation step, we generate m-1 additional datasets which are obtained by increasing the  
21 level of measurement error by  $(1 + \zeta)\sigma_U^2$ , with  $\zeta \in [0; \zeta_m]$  and known. By using OLS, we  
22 would consistently estimate  $\beta_1 \sigma_X^2 / [\sigma_X^2 + (1 + \zeta)\sigma_U^2]$ . We can now focus our attention to a

1 nonlinear regression problem, in which the independent variable is  $\zeta$  and the dependent  
2 one is  $\hat{\beta}_1$ . Asymptotically, we have:

$$3 \quad E(\hat{\beta}_1|\zeta) = G(\zeta) = \frac{\beta_1\sigma_x^2}{\sigma_x^2+(1+\zeta)\sigma_U^2} \quad (4)$$

4 where  $G(0) = \hat{\beta}_{1,naïve}$  (i.e., the slope estimated by OLS) while  $G(-1) = \hat{\beta}_{1,SIMEX}$ . Here  $\hat{\beta}$   
5 represents the estimated value of the  $\beta$  coefficient.

6 We can summarize this methodology in the following steps:

- 7 1- We consider  $m$  to be equal to 7, with  $\zeta$  that set to assume the values [0.5; 1.0; 1.5;  
8 2.0; 2.5; 3.0]. Simulate independent errors with variance equal to  $\zeta\sigma_U^2$  and add it to  
9 the measured variable  $W$ . For each  $\zeta$ , we repeat this step 200 times (we have  
10 explored the sensitivity of the results to the number of replicates, without finding  
11 them to be sensitive).
- 12 2- Estimate the slope  $\beta_1$  for each of the new simulated datasets. For each value of  $\zeta$ ,  
13 we have 200 estimates of the slope; for each one of them, take the average value.  
14 Therefore, we have six points (plus the naïve estimate).
- 15 3- Plot  $\beta_1$  as a function of  $\zeta$  and fit a quadratic function (see Section 5.3.2 in Carroll  
16 et al. (2007) for considerations about the extrapolant function). Once parameterized  
17  $G(\zeta)$ , we want to obtain the value of the function for  $\zeta = -1$ , so that  $G(-1)$  is the  
18 targeted  $\beta_{1,SIMEX}$ .

19 It is clear that we need to have information about the statistical properties of the  
20 measurement errors. As mentioned before, we analyze the impacts of mis-specification of  
21 the statistical properties of the measurement errors on the SIMEX methodology in a  
22 simulation framework. Similar to Chowdhury and Sharma (2007), we start by simulating

1 X from a standard uniform distribution. To generated Y, we use the linear model in  
2 equation (1), and consider four possible sets of values for  $\beta_0$  and  $\beta_1$ :

- 3 1.  $\beta_0 = 0$  and  $\beta_1 = 0.3$ ;
- 4 2.  $\beta_0 = 0$  and  $\beta_1 = 0.8$ ;
- 5 3.  $\beta_0 = 0$  and  $\beta_1 = 2.0$ ;
- 6 4.  $\beta_0 = 0.8$  and  $\beta_1 = 0.8$ .

7 These values could be representative of changes in hydroclimatic variables of annual  
8 to decadal to centennial time scales, including precipitation, temperature, and discharge.  
9 Moreover, we consider  $\varepsilon$  to have a normal distribution with mean equal to zero and standard  
10 deviation equal to 0.1. As far as the measurement error U is concerned, we assume that it  
11 is normally distributed with mean equal to zero and  $\sigma_u = 0.3$ . We perform our analysis  
12 using the SIMEX package (Lederer and Seibold 2019) in R (R Core Team 2022).

13 For case 2 ( $\beta_0 = 0$  and  $\beta_1 = 0.8$ ), we show how SIMEX can correct for the bias in the  
14 regression coefficient in Figure 1. By neglecting the measurement uncertainties, OLS  
15 would return a slope value  $\hat{\beta}_{1,\text{naïve}}$  equal to 0.64. On the other hand, using the SIMEX  
16 approach we have that  $\hat{\beta}_{1,\text{SIMEX}}$  is equal to 0.78. Therefore, given an additive error model  
17 and knowing the measurement error mean, variance, and distribution, SIMEX can remove  
18 the bias from the estimation of the parameters. Similar improvements were observed for  
19 the other three scenarios as well: the SIMEX methodology was able to provide an estimate  
20 of the slope and intercept that was close to the target values, much more so than by using  
21 OLS.

22 We apply SIMEX to seasonal forecasting to explore its potential applicability to  
23 hydroclimatological variables. Figure 2 shows an example related to the prediction of the



1 Southern Oscillation Index (SOI) in February using the sea surface temperature (SST)  
2 averaged over the Niño3.4 region and forecasted from the beginning of February (i.e., 0.5-  
3 lead time forecast) to the beginning of November of the previous year (i.e., 3.5-lead time  
4 forecast). Results are based on the average of 12 members of the GFDL-CM2p5-FLOR-  
5 B01 by the Geophysical Fluid Dynamics Laboratory and part of the North American Multi-  
6 Model Ensemble (NMME; Kirtman et al. 2014). The study period ranges from 1981 to  
7 2020, and we use 1981-2204 for calibration, and 2005-2020 for validation. As a first step,  
8 we compute  $\sigma_u$  from the 12 members of the GFDL model, and we then apply the SIMEX  
9 methodology to estimate the slope and intercept. For short lead times, the regression lines  
10 between SOI and Niño3.4 SST based on SIMEX and OLS are almost identical. As the lead  
11 time increases, there is more separation between the two lines, with a slightly better  
12 performance by SIMEX in terms of root mean squared error. These results point to the  
13 potential suitability of SIMEX for seasonal forecasting, especially for long lead times.

### 14 **3. Results**

15 Let us start by investigating the impact of mis-specification of the mean and variance  
16 of the measurement errors under the above mentioned four scenarios. To accomplish this  
17 task, we apply the SIMEX methodology assuming that  $U$  is normally distributed with mean  
18 equal to zero and standard deviation equal to 0.3, and investigate what happens if we  
19 increase the value of  $\mu_U$  from 0 to 1, and at the same time we have  $\sigma_U$  ranging from 0.3 to  
20 0.6.

21 In Figures 3-6, we have plotted the results for our analysis. On the x-axis we have  
22 increasing values of  $\mu_U$  from 0 to 1, while on the y-axis we have plotted the uncertainties

1 associated with  $\sigma_U$  in percentage of its true value. As far as the intercept is concerned (left  
2 panels), we show  $\Delta_{\beta_0,i} = (\beta_0 - \hat{\beta}_{0,i})/\beta_1$ , where  $\hat{\beta}_{0,i}$  is the estimated intercept using OLS  
3 or SIMEX (indexed by  $i$ ),  $\beta_0$  is its true value, and  $\beta_1$  is value of the true slope. As far as  
4 the slope is concerned (right panels), we show  $\Delta_{\beta_1,i} = 100(\beta_1 - \hat{\beta}_{1,i})/\beta_1$ , where  $\hat{\beta}_{1,i}$  is the  
5 value of the slope estimated from the data using OLS or SIMEX (indexed by  $i$ ). Here we  
6 normalize by the differences in the intercept with respect to the true slope rather than the  
7 true intercept because the latter is set to zero in some of our simulations, affecting the  
8 computation of the ratio. More generally, the normalization of the results with respect to  
9  $\beta_1$  provides indications on the generalization of these results.

10 As expected, the slope is sensitive only to mis-specifications of the standard deviation  
11 of the measurement errors (Figures 3-6). This statement is valid for both OLS and SIMEX  
12 methodologies. We have increasing errors for increasing uncertainties in  $\sigma_U$ , since we only  
13 partially account for the measurement errors: based on equation (3), the reliability ratio  $\lambda$   
14 is still smaller than 1, resulting in an attenuated slope. Overall, with OLS we make an error  
15 between 20% and 50% of the true slope value for uncertainties in  $\sigma_U$  from 0% to 100%.  
16 On the other hand, using SIMEX we have errors in the estimation of the slope ranging from  
17 5% to 30%. Based on our results, the improvements from the use of the SIMEX  
18 methodology with respect to OLS are larger as the uncertainties in  $\sigma_U$  are smaller. On the  
19 other hand, as the uncertainties in the measurement error standard deviation increase, the  
20 advantage of using SIMEX with respect to OLS decreases. This is made clearer from the  
21 bottom-right panels in Figures 3-6, in which the ratio between the OLS and SIMEX errors  
22  $(\Delta_{\beta_1,naïve}/\Delta_{\beta_1,SIMEX})$  decreases from a value greater than 10 for small errors in  $\sigma_U$  to  
23 values closer to 1 (i.e., they perform in a similar way) for large uncertainties in the

1 measurement error. Comparing the results for scenarios 1 to 3 (Figures 3-5), we notice a  
2 slightly reduced sensitivity of both OLS and SIMEX to mis-specification of  $\sigma_U$ . Comparing  
3 Figures 4 and 6, when estimating the slope the presence of a larger intercept tends to  
4 attenuate the sensitivity of both of the approaches to mis-specifications.

5 While the slope is sensitive exclusively to erroneous specification of the measurement  
6 error standard deviation, the intercept is mostly sensitive to mis-specification of its mean  
7  $\mu_U$ . This feature was expected, because errors with a mean different from zero would  
8 introduce an offset that would be compensated by creating or shifting an intercept. For  
9 small values of  $\mu_U$ , the intercept is not sensitive to the error standard deviation. On the  
10 other hand, as  $\mu_U$  increases, the sensitivity of the intercept to increasing uncertainties in  $\sigma_U$   
11 increases: the attenuation in the estimation of the slope tends to reduce the impact of the  
12 increasing values of  $\mu_U$ . Based on Figures 3-6 (bottom-left panel) where we show  
13  $\Delta_{\beta_{0,\text{naïve}}}/\Delta_{\beta_{0,\text{SIMEX}}}$ , the intercept estimated with SIMEX tends to be very close to what  
14 estimated from OLS.

15 Finally, we have looked at the impact of mis-specification of the measurement error  
16 distribution and the role of sample size (Figure 7). Throughout this study, we have assumed  
17 to know the error distribution, and, in particular, that it can be described by a Gaussian  
18 distribution. However, in many hydrometeorological applications we do not have this type  
19 of information. Therefore, we want to consider the case in which we erroneously assume  
20 that the errors are normally distributed, even though they are actually described by another  
21 distribution. In particular, we consider  $U$  to follow Laplace, logistic, lognormal, and  
22 gamma distributions (e.g., Johnson et al. 1994). The Laplace and logistic distributions have  
23 a mean and standard deviation equal to 0 and 0.3, respectively. As far as the lognormal

1 distribution is concerned, these values are also the values of the mean and standard  
2 deviation of the corresponding Gaussian distribution. Finally, the gamma distribution has  
3 a mean of 0.5 and a standard deviation of 0.3. To focus on the impact of higher moment  
4 orders (e.g., skewness, kurtosis), we have linearly transformed the gamma and lognormally  
5 distributed errors to match the values 0.0 and 0.3 for mean and standard deviation,  
6 preserving the shape of their distributions. We consider six sample sizes (i.e.,  $n = 25, 50,$   
7  $100, 200, 400,$  and  $800$ ), and present our results in Figure 7, where the true values of  
8 intercept and slope are  $\beta_0 = 0$  and  $\beta_1 = 0.8$ .

9       Regardless of the distribution and the sample size, the SIMEX method outperforms  
10 OLS in estimating  $\beta_1$ . If we use the median of the  $\hat{\beta}_1$  by SIMEX and OLS as reference, the  
11 gap between the two approaches tends to remain constant, with OLS estimating a slope  
12 value around 0.6, while SIMEX is closer to the target value of 0.8. As we increase the  
13 sample size, the variability associated with the estimation of the slope decreases, especially  
14 as we go from 25 to 100. For sample sizes larger than 100, the marginal improvement in  
15 terms of performance decreases. Neither of the approaches displays a strong dependence  
16 on the distribution of  $U$ . Not surprisingly, when the measurement error follows a Gaussian  
17 distribution, SIMEX performs well. However, even when  $U$  follows a Laplace, gamma or  
18 logistic distribution, the sensitivity of the results is very limited and indeed almost  
19 indistinguishable from what observed for the Gaussian case. The largest departure is for  
20 the lognormal distribution, whose mis-specification has a detectable signal for both SIMEX  
21 and OLS. Therefore, based on these results and in agreement with Carroll et al. (2007), the  
22 SIMEX methodology is robust to violations to the assumption of Gaussian distribution of  
23 the measurement errors.

## 1 **4. Discussion and Conclusions**

2 In studies investigating the existence of a relation between two variables, the presence  
3 of errors in a predictor could significantly affect the results of the analyses. A possible  
4 approach to account for measurement errors is represented by the SIMEX methodology. It  
5 is a simple albeit powerful methodology in the case in which we have information about  
6 the statistical properties of the measurement errors. However, its sensitivity to an erroneous  
7 specification of these errors has received little attention and was the topic of this study. Our  
8 findings can be summarized as follows:

- 9 1- The slope is sensitive to mis-specifications of the measurement error standard  
10 deviation  $\sigma_U$  and insensitive to mis-specification of the measurement error mean  
11  $\mu_U$ .
- 12 2- For low values of the mean  $\mu_U$  (i.e., measurements that are unbiased or have small  
13 biases), the intercept is sensitive only to mis-specifications of the measurement  
14 error mean. As the values of  $\mu_U$  increases, it tends to depend on  $\sigma_U$  as well.
- 15 3- Departure of the statistical distribution of the measurement error from a Gaussian  
16 distribution was found to be not very sensitive to the case of Laplace, logistic and  
17 gamma distribution, while a comparatively worse performance was exhibited by  
18 for the case of the lognormal distribution. Therefore, mis-specifications of the  
19 distribution of the measurement errors do not affect the results of the SIMEX  
20 methodology.

21 In this study we have focused on linear regression. However, SIMEX has been  
22 successfully applied to other and more complex models (e.g., Carroll et al. 2007), and the

1 sensitivity of SIMEX to error mis-specification in these setups should be evaluated in  
2 future studies.

3 Even under mis-specifications of the statistical properties of the measurement errors,  
4 the SIMEX method performs better than the standard OLS. Therefore, our findings suggest  
5 that SIMEX could be a very valuable approach in those applications where the  
6 measurements of a predictor are affected by errors for which limited information about  
7 their statistical characteristics are available, including for sub-seasonal to seasonal  
8 forecasting of hydroclimatic variables (e.g., precipitation, temperature, large-scale climate  
9 indices) or to account for the role of representativeness errors in rain gauges in ground  
10 validation of remote sensing estimates (e.g., Villarini and Krajewski 2008).

11

## 12 **Author Declarations**

13 **Funding:** Gabriele Villarini acknowledges support from the USACE Institute for Water  
14 Resources. Leonardo V. Noto acknowledges financial support for Dario Treppiedi  
15 provided by Consorzio Interuniversitario per l'Idrologia and by Autorità di bacino del  
16 Distretto idrografico della Sicilia.

17

18 **Conflict of interest/Competing interests:** None

19

20 **Ethics approval/declarations:** Not applicable

21

22 **Consent to participate:** Not applicable

23

24 **Consent for publication:** Not applicable

25

26 **Availability of data and material/ Data availability:** Not applicable

27

28 **Code availability:** The codes are available on Github  
29 (<https://github.com/dtreppiedi/simex-mis-specification>).

30

31 **Authors' contributions:** GV. conceptualized the study, and D.T. performed the  
32 simulations and prepared the figures. All authors interpreted the results and wrote the  
33 manuscript.

## References

- Alexeeff, S. E., R. J. Carroll and B. Coull, Spatial measurement error and correction by spatial SIMEX in linear regression models when using predicted air pollution exposures, *Biostatistics*, 17(2), 377-389, 2016.
- Brown, B. W. and R. S. Mariano. Stochastic simulations for inference in nonlinear errors-in-variables models, *Handbook of Statistics*. New York, North Holland, 11, 611-627, 1993..
- Carroll, R. J., D. Ruppert, L. A. Stefanski and C. M. Crainiceanu. *Measurement Error in Nonlinear Models - A Modern Perspective*, CRC Press, 2007..
- Chowdhury, S. and A. Sharma, Mitigating parameter bias in hydrological modelling due to uncertainty in covariates, *Journal of Hydrology*, 340(3-4), 197-204, 2007.
- Chowdhury, S. and A. Sharma, A simulation based approach for representation of rainfall uncertainty in conceptual rainfall runoff models, *Hydrological Research Letters*, 25-8, 2008.
- Cook, J. R. and L. A. Stefanski, Simulation-extrapolation estimation in parametric measurement error models, *Journal of the American Statistical Association*, 89(428), 1314-1328, 1994.
- Devanarayan, V. and L. A. Stefanski, Empirical simulation extrapolation for measurement error models with replicate measurements, *Statistics & Probability Letters*, 59(3), 219-225, 2002.
- Fuller, W. A. *Measurement Error Models*, New York, John Wiley & Sons, 1987..
- Gleser, L. J. Improvements of the naive approach to estimation in nonlinear errors-in-variables regression models, *Statistical Analysis of Error Measurement Models and Application*. P. J. Brown and W. A. Fuller, Providence, Rhode Island, American Mathematics Society, 1990..
- Guolo, A., The SIMEX approach to measurement error correction in meta-analysis with baseline risk as covariate, *Statistics in Medicine*, 33(12), 2062-2076, 2014.
- Hasan, M. M., A. Sharma, F. Johnson, G. Mariethoz and A. Seed, Correcting bias in radar Z-R relationships due to uncertainty in point rain gauge networks, *Journal of Hydrology*, 5191668-1676, 2014.
- Johnson, N. L., S. Kotz and N. Balakrishnan. *Continuous Univariate Distributions*, New York, John Wiley & Sons, 1994..
- Kangas, A. S., Effect of errors-in-variables on coefficients of a growth model and on prediction of growth, *Forest Ecology and Management*, 102(2-3), 203-212, 1998.
- Kinane, S. M., C. R. Montes, T. J. Albaugh and D. R. Mishra, A Model to Estimate Leaf Area Index in Loblolly Pine Plantations Using Landsat 5 and 7 Images, *Remote Sensing*, 13(6), 1140, 2021.
- Kirtman, B. P., D. Min, J. M. Infanti, J. L. Kinter, D. A. Paolino, Q. Zhang, H. van den Dool, S. Saha, M. P. Mendez, E. Becker, P. T. Peng, P. Tripp, J. Huang, D. G. DeWitt, M. K. Tippett, A. G. Barnston, S. H. Li, A. Rosati, S. D. Schubert, M. Rienecker, M. Suarez, Z. E. Li, J. Marshak, Y. K. Lim, J. Tribbia, K. Pegion, W. J. Merryfield, B. Denis and E. F. Wood, The North American Multimodel Ensemble Phase-1 seasonal-to-interannual prediction; Phase-2 toward developing intraseasonal prediction, *Bulletin of the American Meteorological Society*, 95(4), 585-601, 2014.

1           Lauzon, C. B., C. Crainiceanu, B. C. Caffo and B. A. Landman, Assessment of bias  
2 in experimentally measured diffusion tensor imaging parameters using SIMEX, *Magnetic*  
3 *Resonance in Medicine*, 69(3), 891-902, 2013.

4           Lederer, W. and H. Seibold, *simex: SIMEX- And MCSIMEX-Algorithm for*  
5 *Measurement Error Models*, 2019..

6           Oh, E. J., B. E. Shepherd, T. Lumley and P. A. Shaw, Considerations for analysis  
7 of time-to-event outcomes measured with error: Bias and correction with SIMEX, *Statistics*  
8 *in Medicine*, 37(8), 1276-1289, 2018.

9           Ponzi, E., L. F. Keller and S. Muff, The simulation extrapolation technique meets  
10 ecology and evolution: A general and intuitive method to account for measurement error,  
11 *Methods in Ecology and Evolution*, 10(10), 1734-1748, 2019.

12           R Core Team, *R: A Language and Environment for Statistical Computing*, Vienna,  
13 Austria, 2022..

14           Stoklosa, J., C. Daly, S. D. Foster, M. B. Ashcroft and D. I. Warton, A climate of  
15 uncertainty: accounting for error in climate variables for species distribution models,  
16 *Methods in Ecology and Evolution*, 6(4), 412-423, 2015.

17           Villarini, G. and W. F. Krajewski, Empirically-based modeling of spatial sampling  
18 uncertainties associated with rainfall measurements by rain gauges, *Advances in Water*  
19 *Resources*, 31(7), 1015-1023, 2008.

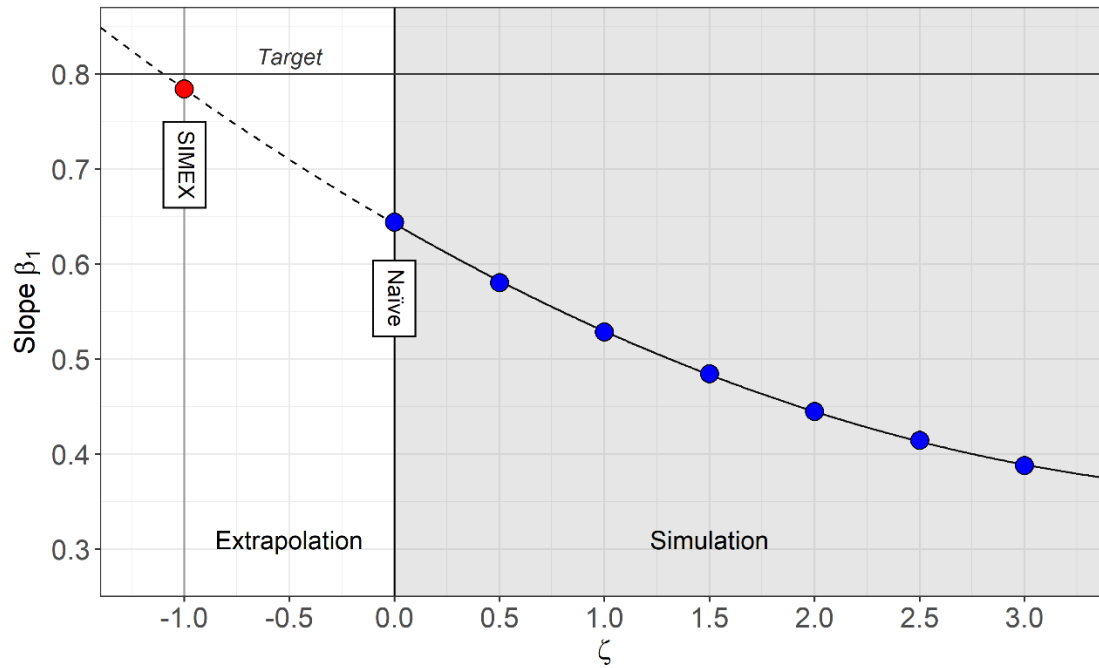
20           Villarini, G., P. V. Mandapaka, W. F. Krajewski and R. J. Moore, Rainfall and  
21 sampling uncertainties: A rain gauge perspective, *Journal of Geophysical Research-*  
22 *Atmospheres*, 113(D11), 2008.

23           Woldemeskel, F. M., A. Sharma, B. Sivakumar and R. Mehrotra, A framework to  
24 quantify GCM uncertainties for use in impact assessment studies, *Journal of Hydrology*,  
25 5191453-1465, 2014.

26

27

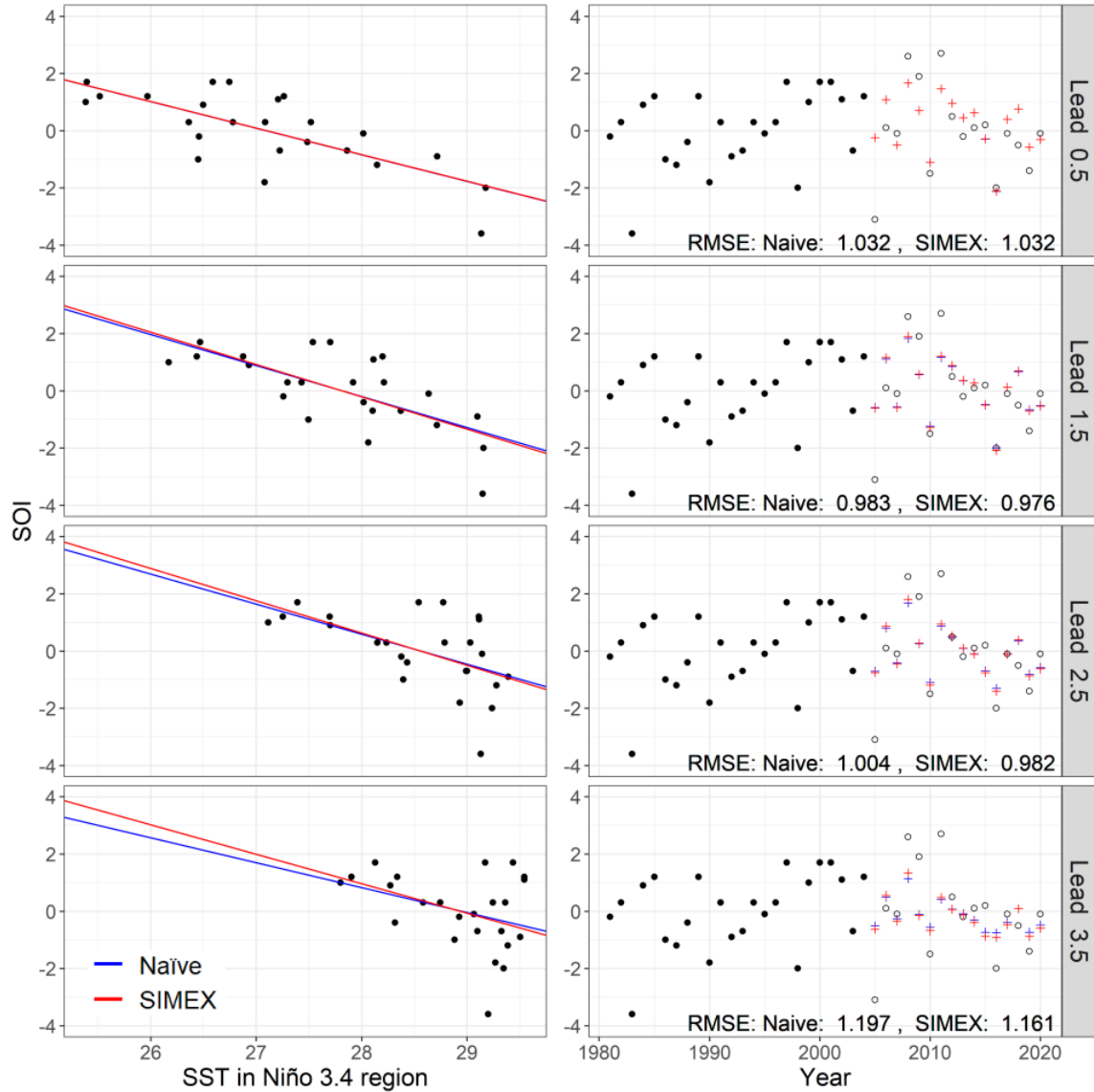




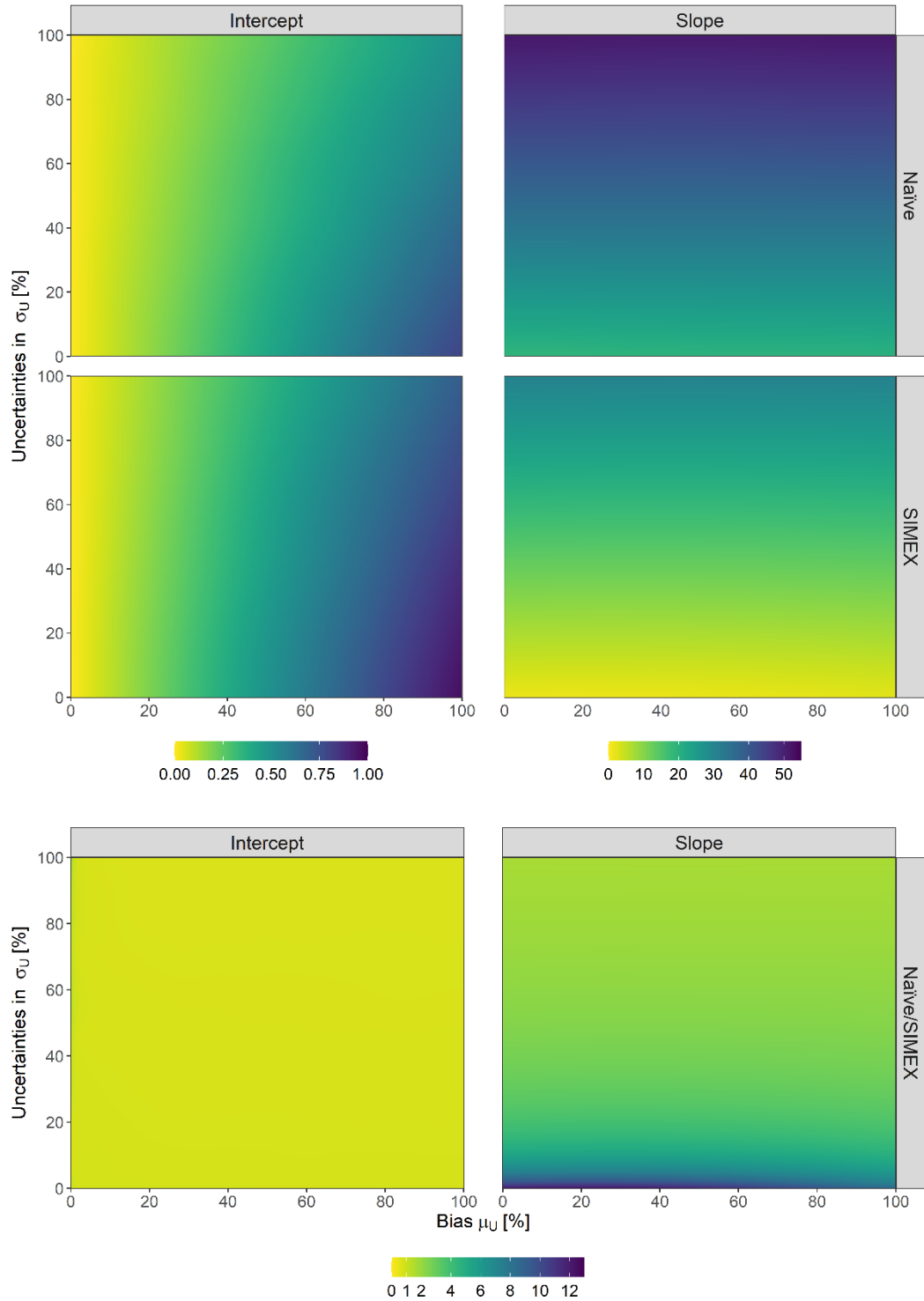
1

2 Figure 1. Plot of the slope  $\beta_1$  as a function of  $\zeta$  for  $\beta_0 = 0$  and  $\beta_1 = 0.8$ . A value of  $\zeta$  equal  
 3 to 0 corresponds to the naïve estimator ( $\hat{\beta}_{1,\text{naïve}} = 0.64$ ), while a value of -1 represents the  
 4 SIMEX estimator ( $\hat{\beta}_{1,\text{SIMEX}} = 0.78$ ).

5

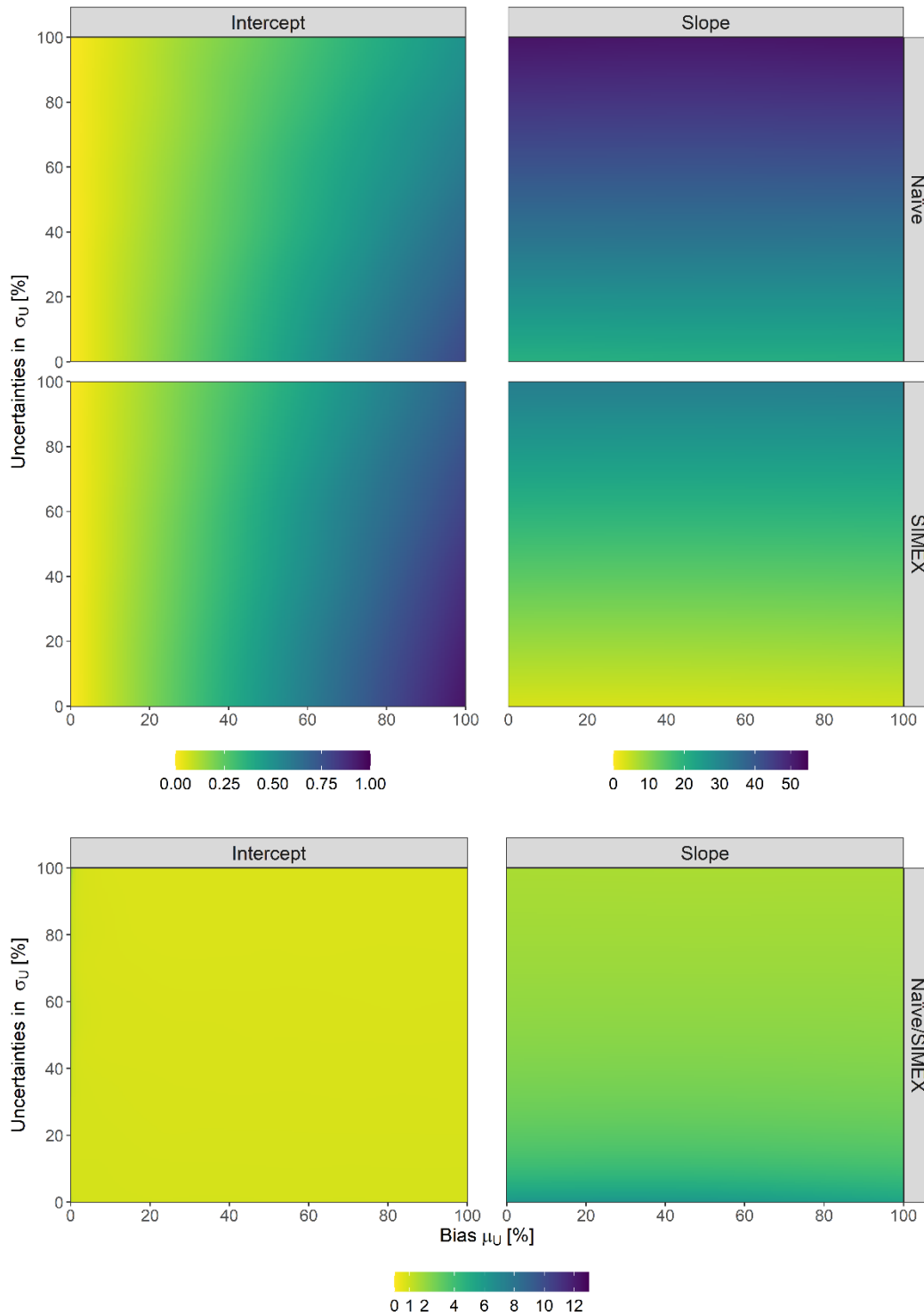


1  
2 Figure 2. Example of the applicability of SIMEX in seasonal forecasting of the Southern  
3 Oscillation Index (SOI) in February as a function of sea surface temperature (SST)  
4 forecasts in the Niño3.4 region. Each row has a different lead time (from shortest to longest  
5 moving from the top to the bottom). The left column shows the scatterplot between the two  
6 variables, together with the regression lines based on OLS (blue) and SIMEX (red). The  
7 right panels show the time series of SOI used for the training of the model (black circles),  
8 the SOI values used for validation (white circles) and the forecasts based on the regression  
9 lines in the left panels, together with the corresponding root mean squared error (RMSE)  
10 values.  
11



1  
2 Figure 3. Sensitivity of the intercept (left panels;  $\Delta\beta_{0,\text{naive}}$  and  $\Delta\beta_{0,\text{SIMEX}}$ , respectively) and  
3 slope (right panels;  $\Delta\beta_{1,\text{naive}}$  and  $\Delta\beta_{1,\text{SIMEX}}$ , respectively) estimates based on the OLS (top  
4 row) and SIMEX (middle row) methodologies to mis-specification of the mean and  
5 standard deviation of the measurement error  $U$ . The value of the true slope  $\beta_1$  is 0.3, while

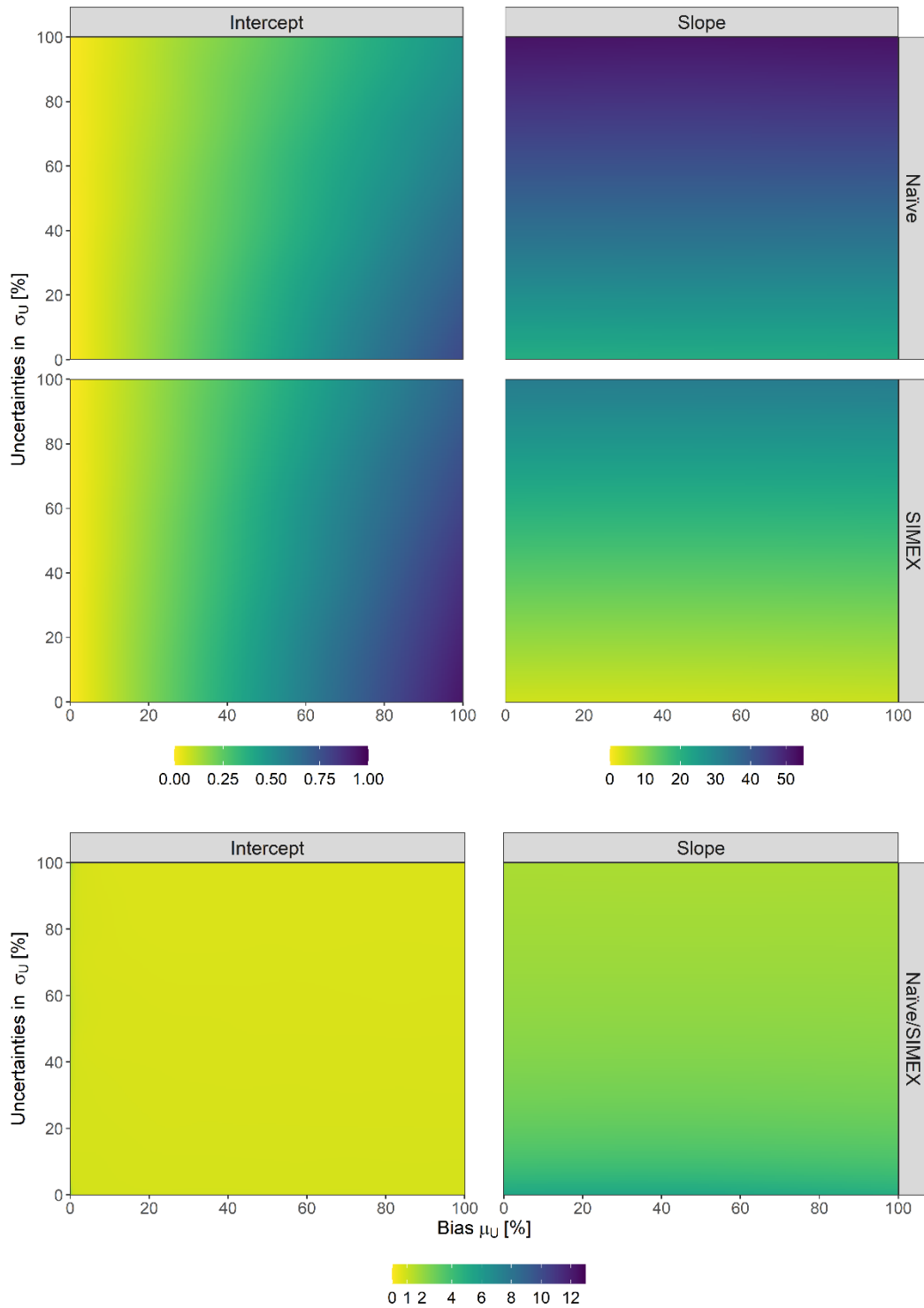
- 1 the value of the true intercept  $\beta_0$  is equal to 0. The panels in the bottom row show
- 2  $\Delta_{\beta_0, \text{naïve}}/\Delta_{\beta_0, \text{SIMEX}}$  (bottom-left) and  $\Delta_{\beta_1, \text{naïve}}/\Delta_{\beta_1, \text{SIMEX}}$  (bottom-right), which represent
- 3 the ratios of the two panels above.



1

2 Figure 4. Same as Figure 3, but for case in which the true slope  $\beta_1$  is 0.8, while the value  
 3 of the true intercept  $\beta_0$  is equal to 0.

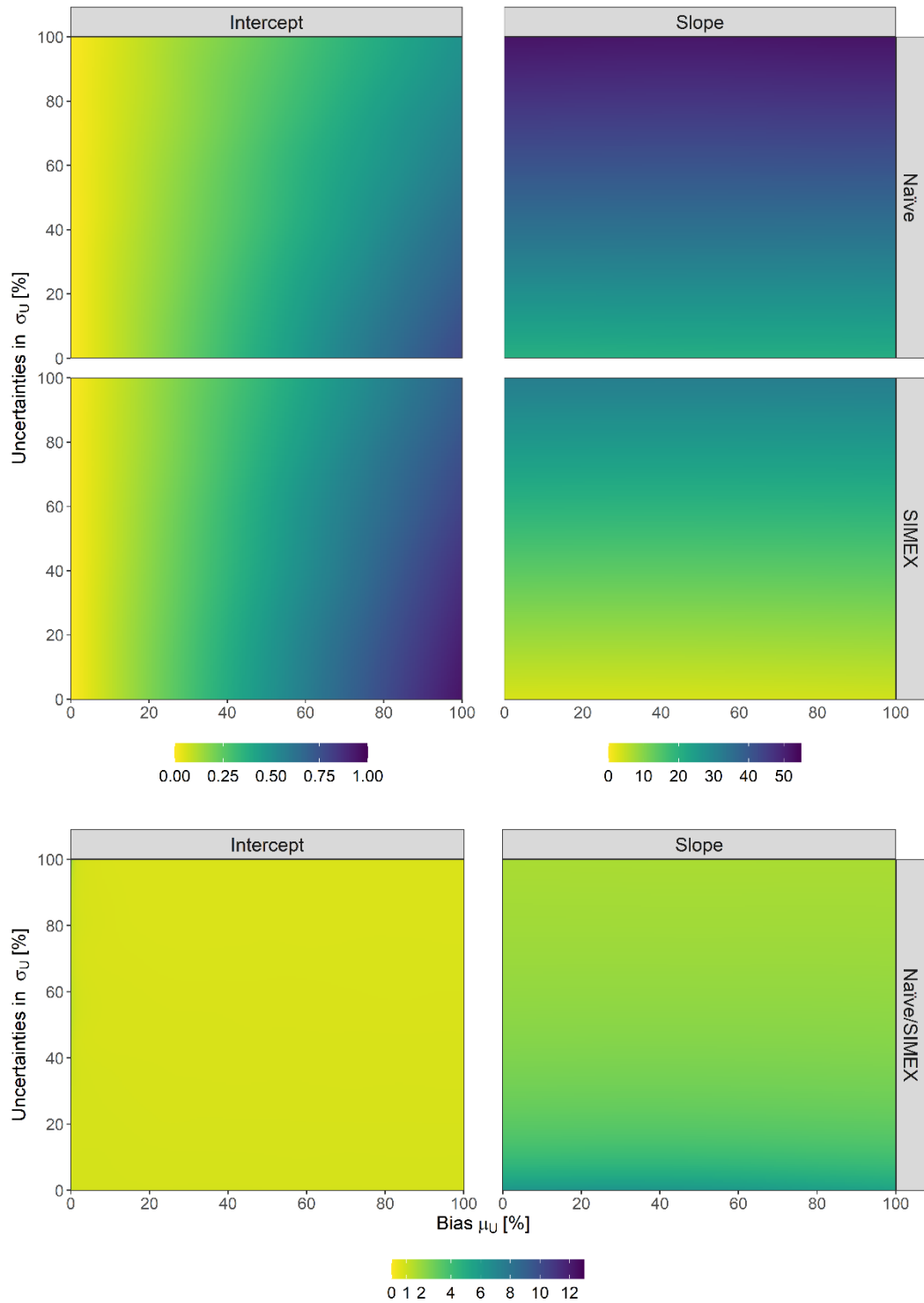
4



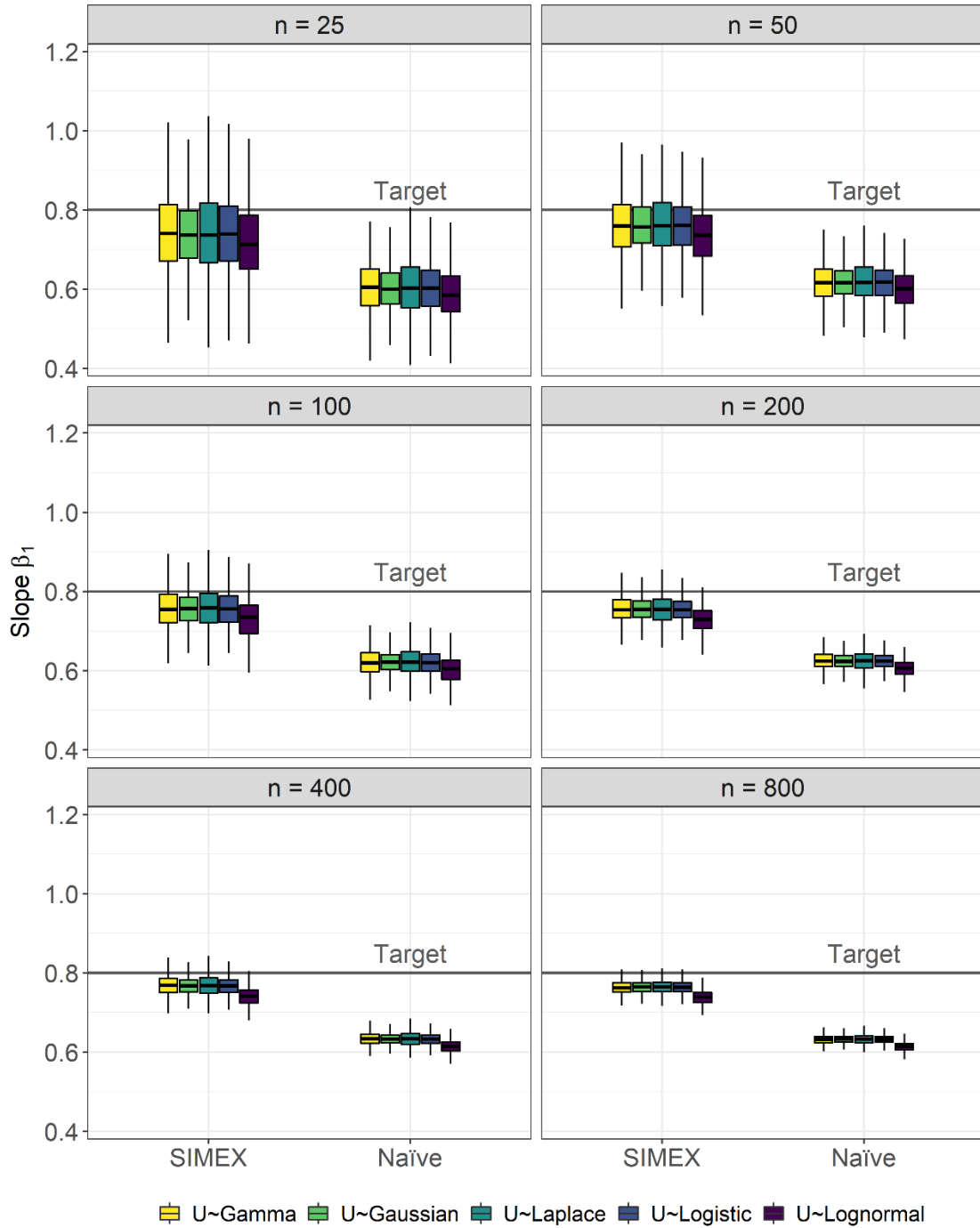
1

2 Figure 5. Same as Figure 3, but for case in which the true slope  $\beta_1$  is 2.0, while the value  
 3 of the true intercept  $\beta_0$  is equal to 0.

4



1  
 2 Figure 6. Same as Figure 3, but for case in which the true slope  $\beta_1$  is 0.8, while the value  
 3 of the true intercept  $\beta_0$  is equal to 0.8.



1

2 Figure 7. Boxplots showing the performance of SIMEX and the naïve estimation of the  
 3 slope  $\beta_1$  as a function of sample size and for different error distributions. In each boxplot,  
 4 the solid line within the box represents the median, while the limits of the box the 25<sup>th</sup> and  
 5 75<sup>th</sup> percentiles; the limits of the whiskers indicate the 5<sup>th</sup> and 95<sup>th</sup> percentiles.

6