

Simultaneous use of observations and deterministic model outputs to forecast persistent stochastic processes

H. Tyralis and D. Koutsoyiannis, Department of Water Resources and Environmental Engineering, National Technical University of Athens (itia.ntua.gr/1401)

1. Abstract

We combine a time series of a geophysical process with the output of a deterministic model, which simulates the aforementioned process in the past also providing future predictions. The purpose is to convert the single prediction of the deterministic model for the future evolution of the time series into a stochastic prediction. The time series is modelled by a stationary persistent normal stochastic process. The output of the deterministic model comprises of the simulation historical part of the process and its deterministic future prediction. The complexity of the deterministic model is assumed to be irrelevant to our framework. A multivariate stochastic process, whose first variable is the true (observable) process and the second variable is a process representing the deterministic model, is formed. The covariance matrix function is computed and the distribution of the unobserved part of the stochastic process is calculated conditional on the observations and the output of the deterministic model.

2. Definitions

We assume that $\{x_{1t}\}, \{x_{2t}\}$, $t = 1, 2, \dots$ are two Hurst-Kolmogorov stochastic processes (HKP) with means μ_1, μ_2 , standard deviations σ_1, σ_2 , autocovariance functions γ_{1k}, γ_{2k} , and autocorrelation functions (ACF) $\rho_{1k} := \gamma_{1k} / \sigma_1, \rho_{2k} := \gamma_{2k} / \sigma_2$, ($k = 0, \pm 1, \pm 2, \dots$).

Then the normal bivariate process $\{x_t = (x_{1t}, x_{2t})^T\}$, $t = 1, 2, \dots$ is a well-balanced HKP if (Amblard et al. 2012)

$$\gamma_{ij}(k) := \text{Cov}\{x_{it}, x_{jt+k}\} = (1/2) \sigma_i \sigma_j (w_{ij}(k-1) - 2w_{ij}(k) + w_{ij}(k+1)) \text{ and } w_{ij}(k) := \rho_{ij} |k|^{H_i+H_j}, \rho_{ij} = 1, \quad (1)$$
$$\rho_{ij} = \rho_{ji} = \rho, (i,j) \in \{1,2\}, \{1,2\}$$

under the restriction $\rho^2 \leq \frac{\Gamma(2H_1+1) \Gamma(2H_2+1) \sin(\pi H_1) \sin(\pi H_2)}{\Gamma^2(H_1+H_2+1) \sin(\pi(H_1+H_2)/2)}$.

The problem of finding and assessing the maximum likelihood estimator for the parameters of the HKP was studied by Tyralis and Koutsoyiannis (2011). The solution of the same problem for the bivariate HKP is more complicated. We assume that there is a record of n observations $x_{1:1:n} := (x_{11} \dots x_{1n})^T$ and $x_{2:1:n} := (x_{21} \dots x_{2n})^T$. The parameters of the bivariate HKP are $\theta = (\mu_1, \mu_2, \sigma_1, \sigma_2, H_1, H_2, \rho)$. We use the terminology of Wei (2006, p.382-427). Hence we have the mean vector $E[x_t] = [\mu_1 \mu_2]^T$ and the lag- k covariance matrix function $\Gamma(k)$:

$$\Gamma(k) := \text{Cov}\{x_t, x_{t+k}\} = \begin{bmatrix} \gamma_{11}(k) & \gamma_{12}(k) \\ \gamma_{21}(k) & \gamma_{22}(k) \end{bmatrix} \quad (2)$$

The covariance matrix of the multivariate normal variable $x_{1:1:n} := [x_1^T \dots x_n^T]^T$ is

$$\Gamma = \begin{bmatrix} \Gamma(0) & \Gamma(1) & \dots & \Gamma(n-1) \\ \Gamma(1) & \Gamma(0) & \dots & \Gamma(n-2) \\ \dots & \dots & \dots & \dots \\ \Gamma(n-1) & \Gamma(n-2) & \dots & \Gamma(0) \end{bmatrix} \quad (3)$$

3. Maximum likelihood estimates

Rearranging the elements of $x_{1:1:n}$ we define the vector $w_{1:1:n} := [x_1^T \dots x_n^T]^T$ with covariance matrix

$$\Sigma = \begin{bmatrix} \Sigma_1 & \Sigma_{12} \\ \Sigma_{21} & \Sigma_2 \end{bmatrix} \quad (4)$$

$$\Sigma_1 := \sigma_1^2 R_{11}, [R_{11}]_{(i,j)} = [R_{11}]_{(j,i)} := \rho_{1(j-i)}, \Sigma_2 := \sigma_2^2 R_{22}, [R_{22}]_{(i,j)} = [R_{22}]_{(j,i)} := \rho_{2(j-i)}, \Sigma_{12} := \rho_{12} \sigma_1 \sigma_2 R_{21}, [R_{21}]_{(i,j)} = [R_{21}]_{(j,i)} := \rho_{21}(j-i) \quad (5)$$

$$\rho_{21}(j-i) := \gamma_{21}(j-i) / (\rho \sigma_1 \sigma_2) = (1/2) (|j-i-1|^{H_1+H_2} - 2|j-i|^{H_1+H_2} + |j-i+1|^{H_1+H_2}) \quad (6)$$

Now we define the vectors

$$e_n = [1 \ 1 \ \dots \ 1]^T, \mu = [\mu_1 e_n^T \ \mu_2 e_n^T]^T \quad (7)$$

The probability distribution function of w is

$$f(w_{1:1:n}) = (2\pi)^{-n} |\Sigma|^{-1/2} \exp\{-(w_{1:1:n} - \mu)^T \Sigma^{-1} (w_{1:1:n} - \mu)\} \quad (8)$$

It is shown that

$$\hat{\sigma}_1 = ((a_1 a_3^2 - \rho a_2 a_1^2) / (n a_3^2))^{1/2}, \hat{\sigma}_2 = ((a_3 a_1^2 - \rho a_2 a_3^2) / (n a_1^2))^{1/2} \quad (9)$$

where

$$a_1 := y_{1:1:n}^T (R_1 - \rho^2 R_{21} R_2^{-1} R_{21}^T)^{-1} y_{1:1:n}, a_2 := y_{2:1:n}^T (R_2 - \rho^2 R_{21} R_1^{-1} R_{21}^T)^{-1} R_{21}^T y_{1:1:n} \quad (10)$$

$$a_3 := y_{2:1:n}^T (R_2 - \rho^2 R_{21} R_1^{-1} R_{21}^T)^{-1} y_{2:1:n} \quad (11)$$

Now substituting (9) in (8) and maximizing the three parameters log-likelihood we obtain $\hat{H}_1, \hat{H}_2, \hat{\rho}$. After substituting these values in (9) we obtain $\hat{\sigma}_1$ and $\hat{\sigma}_2$.

4. Posterior predictive distribution

We assume that $x_{1:1:(n+k)}$ is the output of the deterministic model and $x_{2:1:n}$ is the data observed. We wish to find the distribution of $x_{2:(n+1):(n+m)}$ conditional on $x_{1:1:(n+m)}$ and $x_{2:1:n}$. Assuming that $\{x_t = (x_{1t}, x_{2t})^T\}$, $t = 1, 2, \dots$ is a bivariate HKP, the probability distribution of $w_{1:1:(n+m)}$ is given by (8). The $2(n+m)$ -by- $2(n+m)$ covariance matrix of the process is given by (4) and is partitioned according to (12).

$$\Sigma = \begin{bmatrix} \Sigma_1 & \Sigma_{12} & \Sigma_{122} \\ \Sigma_{21} & \Sigma_{2n} & \Sigma_{2nm} \\ \Sigma_{212} & \Sigma_{2nm} & \Sigma_{2m} \end{bmatrix} = \begin{bmatrix} P_1 & P_{12} \\ P_{21} & P_2 \end{bmatrix} \quad (12)$$

where Σ_{2m} is m -by- m matrix and

$$P_1 = \begin{bmatrix} \Sigma_1 & \Sigma_{121} \\ \Sigma_{211} & \Sigma_{2n} \end{bmatrix}, P_{21} = [\Sigma_{212} \ \Sigma_{2nm}], P_{12} = \begin{bmatrix} \Sigma_{122} \\ \Sigma_{2nm} \end{bmatrix}, P_2 = \Sigma_{2m} \quad (13)$$

Then the posterior predictive distribution of $x_{2:(n+1):(n+m)}$ conditional on $x_{1:1:(n+m)}$, $x_{2:1:n}$ and θ is

$$f(x_{2:(n+1):(n+m)}) = (2\pi\sigma^2)^{-m/2} |R_{m|n}|^{-1/2} \exp\{(-1/2\sigma^2) (x_{2:(n+1):(n+m)} - \mu_{m|n})^T R_{m|n}^{-1} (x_{2:(n+1):(n+m)} - \mu_{m|n})\} \quad (14)$$

$$\mu_{m|n} := \mu_2 e_m + P_{21} P_1^{-1} ([x_{1:1:(n+m)}^T \ x_{2:1:n}^T]^T - [\mu_1 e_{n+m}^T \ \mu_2 e_n^T]^T) \quad (15)$$

$$R_{m|n} := P_2 - P_{21} P_1^{-1} P_{12} \quad (16)$$

Here we mention that in the following θ will be considered known and equal to its maximum likelihood estimate. In a Bayesian setting we would assume that θ is a random variable, but this is out of the scope of this study. In the Bayesian setting the uncertainty of the prediction would increase (see e.g. Tyralis and Koutsoyiannis, 2013a). The variables that will be examined in the following will be considered normal. For truncated normal variables the interested reader is referred to Horrace (2005) and Tyralis and Koutsoyiannis (2013). The examination of non-normal variables is out of the scope of this study as well. For more details on the method and how it is compared to the methods of Krzysztofowicz (1999) and Wang et al. (2009) see Tyralis and Koutsoyiannis (2013b).

5. Methodology

We applied our methodology on global temperature data and precipitation data shown in Table 1. These data are modeled by a Hurst-Kolmogorov process (Koutsoyiannis, 2011). Koutsoyiannis and Montanari, 2007). The deterministic models used in the study were the General Circulation Models (GCMs). We used the 20C3M for the calibration of the model and the SRES scenarios A1B, B1, A2 (Table 2) were taken into consideration for the prediction of the stochastic model. The specific GCMs that were used in the study are shown in Table 3. Tables 4-7 show the maximum likelihood estimates of the bivariate HKP $\{x_t = (x_{1t}, x_{2t})^T\}$, where $\{x_{1t}\}$ is the process which models the GCM and $\{x_{2t}\}$ is the process which models the observations. The time interval for the calibration spans from the maximum starting year of the corresponding 20C3M scenario and the observed data to the minimum of the corresponding 20C3M scenario and the observed data. We also examined the case where the parameters are estimated separately. Specifically the $\{x_{1t}\}$, $\{x_{2t}\}$ are assumed to be univariate HKPs and their parameters are estimated as in Tyralis and Koutsoyiannis (2011). The sample cross-correlation function is used in this case to estimate ρ .

Using the simultaneous maximum likelihood estimate of ρ , we obtain the posterior predictive distribution of $x_{2:(n+1):(n+m)}$ conditional on $x_{1:1:(n+m)}$, $x_{2:1:n}$ and θ from (14). The other parameters of the bivariate process are estimated again assuming that $\{x_{1t}\}$, $\{x_{2t}\}$ are univariate HKPs, however in this case we use the whole sample, starting from the common starting year of $\{x_{1t}\}$ and $\{x_{2t}\}$ until the year 2100 for the $\{x_{1t}\}$ parameter estimates and the common end year of the corresponding 20C3M scenario and $\{x_{2t}\}$ parameter estimates. The samples from the posterior predictive probability of x_{1t}, x_{2t} , $t = n+1, n+2, \dots$, were used to obtain samples for the variable of interest $x_t^{(30)}$ given by (17) following the framework in Koutsoyiannis et al. (2007).

$$x_t^{(30)} := (1/30) (\sum_{l=t-29}^n x_{1l} + \sum_{l=t-n+1}^t x_{2l}), t = n+1, \dots, n+29 \text{ and } x_t^{(30)} := (1/30) (\sum_{l=t-n+30}^t x_{1l} + \sum_{l=t-29}^{t-1} x_{2l}), t = n+30, n+31, \dots \quad (17)$$

6. Examined datasets

Table 1. Study historical time series. Table 2. IPCC scenarios and their relevance to the study.

Table 3. Main characteristics of the GCMs used in the study. Table 4. Maximum likelihood estimates for the parameters of the bivariate HKP for the CRU combined land (CRUTEM) and marine temperature anomalies.

7. Maximum likelihood estimates for temperature datasets parameters

Table 4. Maximum likelihood estimates for the parameters of the bivariate HKP for the CRU combined land (CRUTEM) and marine temperature anomalies. Table 5. Maximum likelihood estimates for the parameters of the bivariate HKP for the NOAA annual global land and ocean temperature anomalies.

Highlighted are the cases whose results were used at Figures 1-12

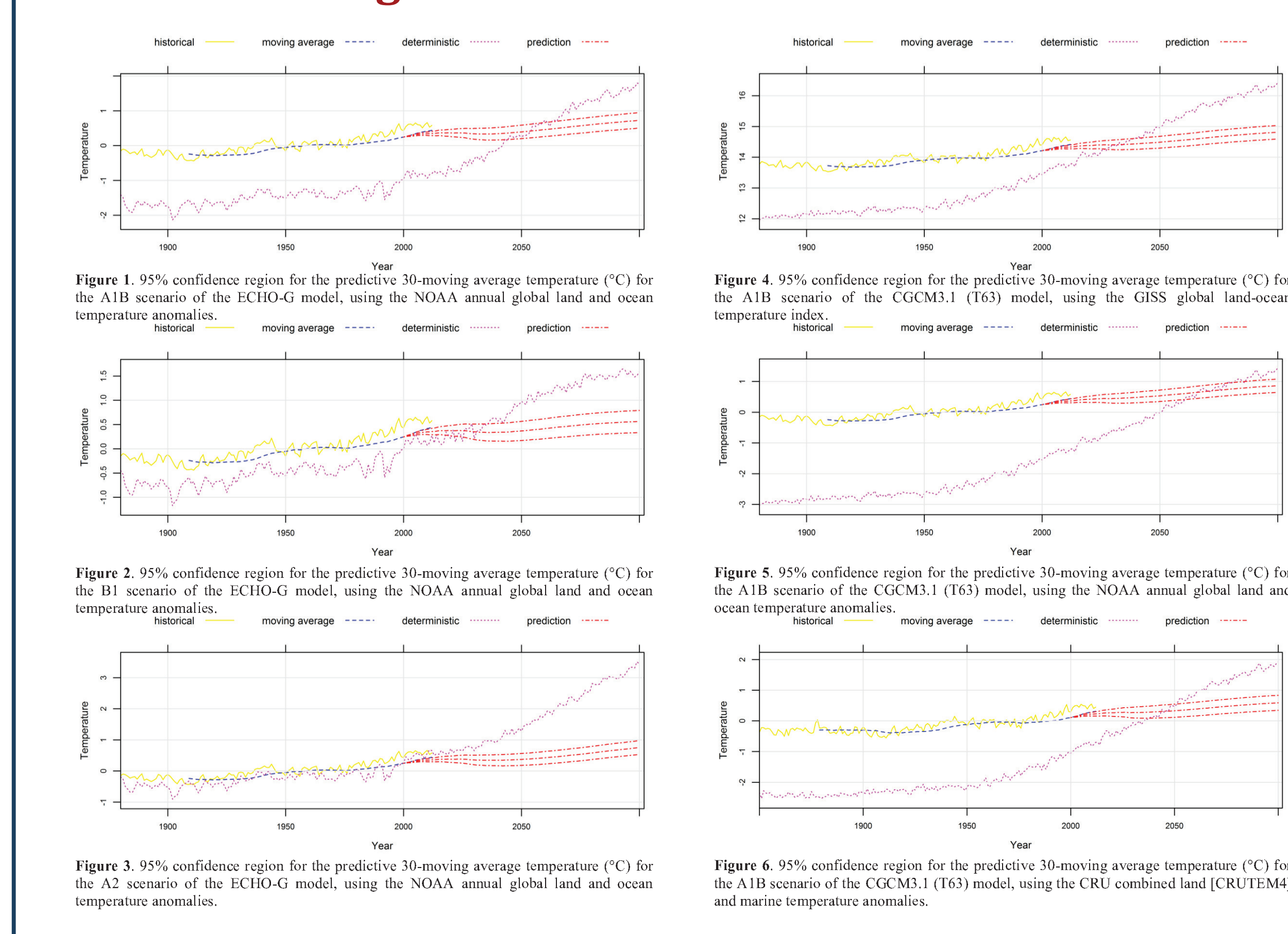
8. Maximum likelihood estimates for the precipitation dataset parameters

Table 7. Maximum likelihood estimates for the parameters of the bivariate HKP for the CRU precipitation over land areas.

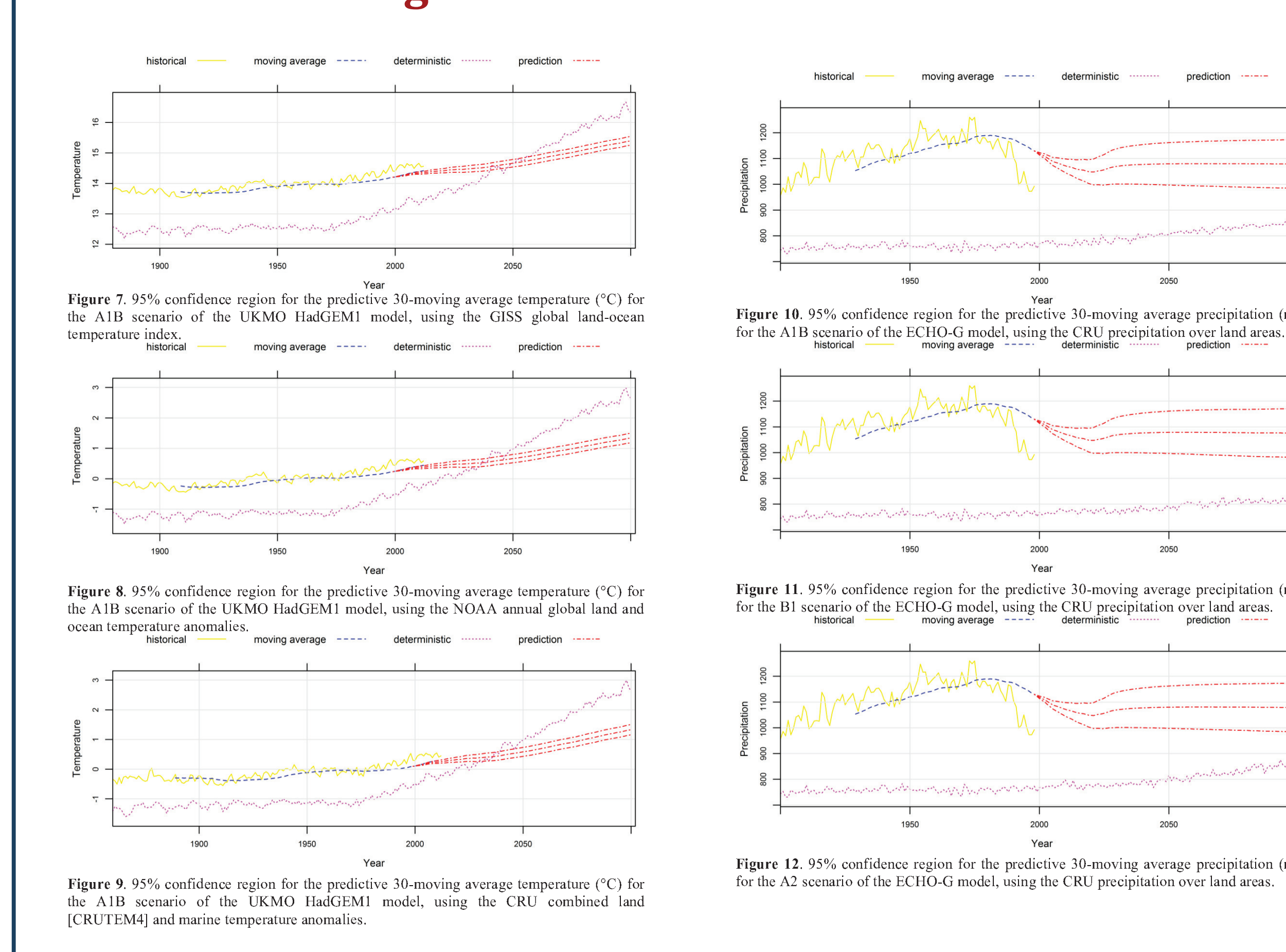
Highlighted are the cases whose results were used at Figures 1-12

Separate MLEs are not used in the study, because the parameters are not orthogonal

9. Confidence regions for future climate



10. Confidence regions for future climate



11. Conclusions

- We derive a new estimator for the parameters of the bivariate HKP.
After modelling the observed datasets and the output of GCMs using the bivariate HKP we estimate the parameters of the process.
Using the estimated values of the parameters we provide stochastic prediction of the future climate combining the projections of the GCMs and their corresponding hindcasts with the observed time series.
The estimated values of the cross-correlation between the historical datasets (at global scale) and the hindcasts of the GCMs range from 0 to 0.4, showing that the information added by the GCMs to that contained in the historical datasets is not substantial.
The upper bound of the 95% confidence region of the climatic value of temperature at year 2100 is estimated to about 1°C more than the current value of this climatic variable.
For the precipitation dataset the estimated value of the cross-correlations between the historical datasets and the hindcasts of the GCMs is almost equal to 0. This means that the output of the GCM has no effect on the stochastic predictions.
We emphasize that the estimation of the stochastic model parameters should better be performed using only data that were not used in the GCM fitting/tuning, i.e. for the period after 2000. This would correspond to the so-called split-sample technique, which avoids possible model overfitting on the available data. However this would increase considerably the uncertainty of the estimators of the parameters of the models and practically would result in total neglect of the GCM predictions. Hence we decided to approach the problem more conservatively.
Our approach is an extension of previous studies, which exploited the outputs of deterministic models combined with historical dataset, on persistent stochastic processes.

12. References

Amblard PO, Couerjolly JF, Lavancier F, Philippe A (2012) Basic properties of the multivariate fractional Brownian motion. arXiv:1007.0828v2
Carter TR, Hulme M, Lal M (1999) Guidelines on the use of scenario data for climate impact and adaptation assessment, task group on scenarios for climate impact assessment. Intergovernmental Panel on Climate Change
Heeger G, Mehl G, Covey C, Latif M, McAvaney B, Stouffer R (2003) 20C3M: CMIP collecting data from 20th century coupled model simulations. Exchanges 26, 9(1), Int. CLIVAR Project Office
Horrace W (2005) Some results on the multivariate truncated normal distribution. Journal of Multivariate Analysis 94(1):209-221. doi:10.1016/j.jmva.2004.10.007
Koutsoyiannis D (2011) Hurst-Kolmogorov dynamics and uncertainty. JAWRA Journal of the American Water Resources Association 47(3):481-495. doi:10.1111/j.1752-1688.2011.00543.x
Koutsoyiannis D, Efstratiadis A, Georgakakos K (2007) Uncertainty Assessment of Future Hydroclimatic Predictions: A Comparison of Probabilistic and Scenario-Based Approaches. Journal of Hydrometeorology 8(3):261-281. doi:10.1175/JHM576.1
Koutsoyiannis D, Montanari A (2007) Statistical analysis of hydroclimatic time series: Uncertainty and insights. Water Resources Research 43(5) W05429. doi:10.1029/2006WR005592
Krzysztofowicz R (1999) Bayesian Forecasting via Deterministic Model. Risk Analysis 19(4):739-749. doi:10.1111/j.1539-6924.1999.tb00443.x
Leggett J, Pepper WJ, Swart RJ (1992) Emissions scenarios for the IPCC: an update. In: Climate Change 1992: The Supplementary Report to the IPCC Scientific Assessment (ed. by J. T. Houghton, B. A. Callander & S. K. Varney), 75-95. Cambridge University Press, Cambridge, UK
Nakicenovic N, Swart R (eds) (1999) IPCC Special Report on Emissions Scenarios. Intergovernmental Panel on Climate Change. Available online at: http://www.grida.no/publications/other/ipcc_sr/?src=/climate/ipcc/emission/
Tyralis H, Koutsoyiannis D (2011) Simultaneous estimation of the parameters of the Hurst-Kolmogorov stochastic process. Stochastic Environmental Research & Risk Assessment 25(1):21-33. doi:10.1007/s00477-010-0408-x
Tyralis H, Koutsoyiannis D (2013a) A Bayesian statistical model for deriving the predictive distribution of hydroclimatic variables. Climate Dynamics. doi:10.1007/s00382-013-1804-y
Tyralis H, Koutsoyiannis D (2013b) On the prediction of persistent processes using the output of deterministic models. In preparation
Wang QJ, Robertson DE, Chiew FHS (2009) A Bayesian joint probability modeling approach for seasonal forecasting of streamflows at multiple sites. Water Resources Research 45(5). doi:10.1029/2008WR007355
Wei WWS (2006) Time Series Analysis, Univariate and Multivariate Methods. second edition, Pearson Addison Wesley, Chichester