

Rothamsted Repository Download

A - Papers appearing in refereed journals

Curceac, S., Atkinson, P. M., Milne, A. E., Wu, L. and Harris, P. 2020. Adjusting for conditional bias in process model simulations of hydrological extremes: an experiment using the North Wyke Farm Platform. *Frontiers in Artificial Intelligence*. 3 (82). <https://doi.org/10.3389/frai.2020.565859>

The publisher's version can be accessed at:

- <https://doi.org/10.3389/frai.2020.565859>
- <https://www.frontiersin.org/articles/10.3389/frai.2020.565859/abstract>

The output can be accessed at:

<https://repository.rothamsted.ac.uk/item/9820q/adjusting-for-conditional-bias-in-process-model-simulations-of-hydrological-extremes-an-experiment-using-the-north-wyke-farm-platform>.

© 9 October 2020, Please contact library@rothamsted.ac.uk for copyright queries.

Adjusting for conditional bias in process model simulations of hydrological extremes: an experiment using the North Wyke Farm Platform

1

2

3 **Stelian Curceac^{1*}, Peter M. Atkinson^{2,3,4}, Alice Milne⁵, Lianhai Wu¹, Paul Harris¹**

4 ¹ Rothamsted Research, Department of Sustainable Agriculture Sciences, North Wyke EX20 2SB, Devon,
5 UK.

6 ²Lancaster Environment Centre, Lancaster University, Bailrigg, Lancaster LA1 4YQ, UK.

7 ³Geography and Environment, University of Southampton, Highfield, Southampton SO17 1BJ, UK.

8 ⁴State Key Laboratory of Resources and Environmental Information System, Institute of Geographical
9 Sciences and Natural Resources Research, Chinese Academy of Sciences, Beijing 100101, China.

10 ⁵Rothamsted Research, Department of Sustainable Agriculture Sciences, Harpenden AL5 2JQ, UK

11

12 ***Correspondence:**

13 Stelian Curceac

14 stelian.curceac@rothamsted.ac.uk

15

16 **Keywords: peak flow, conditional extreme model, extreme learning machine, process-based**
 17 **model, hybrid, grassland agriculture.**

18 **Abstract**

19 Peak flow events can lead to flooding which can have negative impacts on human life and ecosystem
 20 services. Therefore, accurate forecasting of such peak flows is important. Physically-based process
 21 models are commonly used to simulate water flow, but they often under-predict peak events (i.e., are
 22 conditionally biased), undermining their suitability for use in flood forecasting. In this research, we
 23 explored methods to increase the accuracy of peak flow simulations from a process-based model by
 24 combining the model's output with: (a) a semi-parametric conditional extreme model and (b) an
 25 extreme learning machine model. The proposed 3-model hybrid approach was evaluated using fine
 26 temporal resolution water flow data from a sub-catchment of the North Wyke Farm Platform, a
 27 grassland research station in south-west England, UK. The hybrid model was assessed objectively
 28 against its simpler constituent models using a jackknife evaluation procedure with several error and
 29 agreement indices. The proposed hybrid approach was better able to capture the dynamics of the flow
 30 process and, thereby, increase prediction accuracy of the peak flow events.

31 **1 Introduction**

32 In the UK, the estimated yearly cost of damages caused by floods is over £1 billion (Collet et al., 2017).
 33 Accurate and reliable forecasting of extreme flow events is crucial for planning and implementing
 34 measures to mitigate their effects and so protect lives, properties and services. The magnitude and
 35 frequency of floods is likely to increase as a result of climate change (Bates et al., 2008; Field et al.,
 36 2012; Kundzewicz et al., 2007) and this could push ecosystems beyond the threshold of normal
 37 disturbance (Thibault & Brown, 2008). Increased runoff and flooding intensify erosion and result in
 38 higher sediment and nutrient losses that can lead to soil degradation and high concentrations of
 39 pollutants in water courses (Bouraoui et al., 2004).

40 Over recent decades, different approaches have been proposed for more accurate modelling and
 41 forecasting of peak flows with reduced uncertainty. The two main methods of modelling hydrological
 42 variables are physically-based models and statistical models. However, there is an increasing trend
 43 towards combining these approaches in hybrid models. One of the most common ways to do this is to
 44 post-process statistically an ensemble of forecasts from process-based models (e.g., Cloke and
 45 Pappenberger, 2009; Li et al., 2017). Bayesian methods using climate indices (Bradley et al., 2015),
 46 stochastic data-driven methods on wavelet decomposed series (Quilty et al., 2019), Bayesian model
 47 averaging (Raftery et al., 2005), extended logistic regression (Roulin and Vannitsem, 2011), quantile
 48 regression (López López et al., 2014), bias correction (Li et al., 2019) and nearest neighbor resampling
 49 for uncertainty estimation (Sikorska et al., 2015) are among the many post-processing techniques
 50 described in the literature. Examples of combining a process-based model with more than one statistical
 51 or machine learning model can be found in Bogner et al. (2017), Papacharalampous et al. (2019) and
 52 Tyralis et al. (2019). The usefulness of combining deterministic and stochastic models (Box and
 53 Jenkins, 1976) in real-time flood forecasting was reported by Toth et al. (1999), while the performance
 54 of various post-processing techniques according to the level of flow was investigated in Bogner et al.
 55 (2016) and Papacharalampous et al. (2019). Hybrid methods for water flow (streamflow) forecasting
 56 also include the combination of classical statistical methods with more data-driven, machine-learning
 57 methods such as artificial neural networks (ANNs) (Chen et al., 2018; Yaseen et al., 2016; Zhou et al.,
 58 2018), discrete wavelet transforms and support vector machines (Kisi and Cimen, 2011), and coupling

59 ANNs with autoregressive techniques (Fathian et al., 2019). The effect of catchment characteristics on
60 the predictive performance of two different statistical models was discussed in Dogulu et al. (2015).

61 Hydrological process-based models (PBMs) are traditionally used for streamflow modelling and
62 forecasting, where under-prediction of peak flows is a common issue (e.g., Lane et al., 2019;
63 Wijayarathne and Coulibaly, 2020). The PBM performance can suffer from uncertainty due to both
64 random and systematic errors. Both random and systematic errors can arise in the estimated model
65 parameters and measured input variables. However, of particular interest is a type of systematic error
66 (or bias) called conditional bias that depends on flow magnitude. That is, the structure and parameters
67 of the model can generalise the outputs leading to conditional bias, specifically under-prediction of
68 large values and over-prediction of small values; an effect similar in nature to that of having a support
69 that is larger than ideal. Alternatively, data-driven methods may be used, especially when the initial
70 conditions and the parameters of the physical model are difficult to estimate or when the length and/or
71 quality of the data are insufficient for a reliable model calibration.

72 In this research, we explored combining statistical and machine learning techniques with flow
73 simulations obtained from a PBM to increase the accuracy of forecasting peak flow events.
74 Specifically, we considered the semi-parametric, conditional extreme model (CEM) of Heffernan and
75 Tawn (2004) (a statistical model) and the extreme learning machine (ELM) of Huang et al. (2006) (a
76 machine learning model). The proposed approach is considered a generic solution for enhancing any
77 given hydrological PBM.

78 The CEM is appropriate for describing the probability that one or multiple variables are extreme and
79 has been applied widely for flood risk analysis (Mendes and Pericchi, 2009; Lamb et al., 2010; Keef
80 et al., 2013; Zheng et al., 2014). A significant property of the CEM is that it is flexible in modelling
81 different dependence structures, such as the dependence of different variables at the same site or the
82 dependence of the same variable at different sites. A key assumption of the application of the CEM is
83 that the extremes of each variable must be independent and, consequently, cannot be used to model
84 peak flow events that have a duration of several consecutive days and, therefore, exhibit temporal
85 dependence. For this reason, the maximum flow during each event was modelled using the CEM while
86 all other peaks were modelled using the ELM (and, thus, a 3-model rather than a 2-model hybrid is
87 proposed).

88 The ELM model is ANN-based and has been used in various areas of water resources engineering,
89 with a recent focus on water flow (see Yaseen et al., 2019 for an extensive review). In this context, it
90 has been shown to increase accuracy and reduce computational time compared to commonly used
91 benchmark models (Lima et al., 2015) and to other ANN models (Deo and Şahin, 2016).

92 The resultant 3-model hybrid was evaluated empirically using measured flow data from a sub-
93 catchment of the North Wyke Farm Platform, a grassland research facility in south-west England (Orr
94 et al., 2016). To our knowledge, no study to-date has used the CEM and the ELM to improve the
95 simulation of peak flow events obtained from a PBM, or in which they are combined. The proposed
96 methodology builds on the modelled dependence structure between measured and PBM-simulated
97 peak flow events and uses this relationship to obtain a more accurate representation of these events.

98 **2 Methods**

99 This section presents a general description of the CEM (Heffernan and Tawn, 2004) and the ELM
100 (Huang et al., 2006) and explains how they can be applied to peak flow events obtained from a chosen
101 PBM (described in Section 3.2) in a hybrid context. The flow threshold, above which the simulated

102 and the observed data are considered as possible peaks, is determined based on Generalised Pareto
 103 Distribution (GPD) stability plots of the PBM simulated values (Curceac et al., 2020). The performance
 104 of the proposed hybrid approach is evaluated using a jackknife procedure and by calculating several
 105 error and agreement indices.

106 2.1 Generalised Pareto Distribution (GPD)

107 We characterise peak flow events by fitting the GP distribution to the extreme flow above a certain
 108 threshold. The cumulative distribution function (CDF) of the iid excesses over an appropriately high
 109 threshold u for the GPD is:

$$110 \quad G(x) = \Pr(X - u < x | X > u) = \begin{cases} 1 - \left(1 + \frac{\xi(x - u)}{\sigma}\right)^{-\frac{1}{\xi}}, & \xi \neq 0 \\ 1 - e^{-\frac{x-u}{\sigma}}, & \xi = 0 \end{cases}$$

111 where x , for this study, is the peak flow in mm d^{-1} , u is the location parameter, σ is the scale parameter
 112 and ξ is the shape parameter. The value of the shape parameter defines the type of distribution from
 113 the GPD family; that is, $\xi = 0$ refers to the exponential distribution, the distribution has an upper bound
 114 of $u - \sigma/\xi$ when $\xi < 0$ and has no upper limit when $\xi \geq 0$.

115 The first step in modelling the exceedances is to select a threshold over which peaks in flow are
 116 considered extreme. The next step is to ensure that the peaks above it are independent (so as to conform
 117 with iid) and estimate the scale and shape parameters. The selection of the threshold is a crucial step
 118 in GPD extreme value analysis and is basically a trade-off between bias (low threshold-large sample
 119 size) and variance (high threshold-small sample size).

120 The flow threshold in this research was selected based on the simulated flow from the study's PBM
 121 using an automated threshold stability method (Curceac et al., 2020) (Section 2.2) and the same
 122 threshold was used for the measured flow data. The GP model was fitted initially independently to the
 123 simulated and observed peak flows and the conditional dependence structure between them was
 124 estimated using the CEM (Section 2.3).

125 2.2 GPD Threshold Selection

126 If the GPD is an appropriate model for the excesses above a threshold u , then for all larger thresholds
 127 $u^* > u$ it will also be suitable with the shape parameter being relatively constant (Coles, 2001; Scarrott
 128 & MacDonald, 2012). That is, it is the approximately linear and horizontal segment on a plot of shape
 129 parameter against threshold. This does not apply for the scale parameter σ_{u^*} , which changes with the
 130 threshold $\sigma_{u^*} = \sigma_u + \xi(u^* - u)$. However, the modified scale parameter $\sigma_1 = \sigma_{u^*} - \xi u$ remains
 131 relatively constant. Therefore, following Curceac et al. (2020), we fitted a cubic smoothing spline to
 132 this plot and calculated the rate of change at each of m consecutive steps. The cubic smoothing spline
 133 estimate \hat{f} of a function f in the model $Y_i = f(x_i) + \varepsilon_i$, is defined as the minimizer of
 134 $\sum_{i=1}^n \{Y_i - \hat{f}(x_i)\}^2 + \lambda \int \hat{f}''(x)^2 dx$, where λ is the smoothing parameter. The minimum change rate
 135 locates the part of the plot where the shape and the modified scale parameters reach a plateau.

136 **2.3 Conditional Extreme Model (CEM)**

137 For a continuous d -dimensional vector variable $X = (X_1, \dots, X_d)$ with unknown distribution function
 138 $F(x)$, the CEM describes the distribution function of X when it is extreme in at least one component.
 139 In other words, it describes the conditional distribution of $X_{-i}|X_i > u_{X_i}$, where X_{-i} is the vector
 140 variable X without the component X_i .

141 After estimating the marginal distribution of each $X_i, i = 1, \dots, d$ (Section 2.1), and before estimating
 142 the extremal dependence, the variables are transformed so that they follow the same distribution. This
 143 process is called marginal standardization and is used to distinguish the marginal behaviour from the
 144 dependence structure (Drees and Janßen, 2017). The data can be transformed to either Gumbel margins
 145 to describe the positive dependence or to a Laplace marginal distribution which, due to its exponential
 146 tail and symmetry, captures both positive and negative dependence (Keef et al., 2013). The initial
 147 vector variable X is, therefore, transformed as:

$$148 \quad f(x) = \begin{cases} \log\{2F_{X_i}(X_i)\}, & X_i < F_{X_i}^{-1}(0.5) \\ -\log\{2[1 - 2F_{X_i}(X_i)]\}, & X_i \geq F_{X_i}^{-1}(0.5) \end{cases}$$

149 where $F_{X_i}^{-1}$ is the inverse cumulative distribution function of X_i . The resulting vector variable $Y =$
 150 (Y_1, \dots, Y_d) , therefore, has Laplace margins with:

$$151 \quad \Pr(Y_i \leq y) = F_{Y_i}(y) = \begin{cases} \frac{1}{2} \exp(y), & y < 0 \\ 1 - \frac{1}{2} \exp(-y), & y \geq 0 \end{cases}$$

152 The dependence model considers the asymptotics of the conditional distribution $\Pr(Y_{-i} \leq y_{-i} | Y_i =$
 153 $y_i)$, where for $y_i \rightarrow \infty$, the increase of y_{-i} must result in non-degenerate margins. For this, assume the
 154 normalizing functions $a_{|i}(y_i)$ and $b_{|i}(y_i)$, that have the same dimension as Y_{-i} and for which:

$$155 \quad \lim_{y_i \rightarrow \infty} \left[\Pr \left\{ \frac{Y_{-i} - a_{|i}(y_i)}{b_{|i}(y_i)} \leq z_{|i} \mid Y_i = y_i \right\} \right] = G_{|i}(z_{|i})$$

156 where the limit distribution $G_{|i}$ has non-degenerate marginals $G_{j|i}$ for all $j \neq i$. Therefore, the random
 157 variable $Z_{|i} = \frac{Y_{-i} - a_{|i}(y_i)}{b_{|i}(y_i)}$ is independent of $Y_i > u_{Y_i}$ and has distribution function $G_{|i}$. The location
 158 $a_{|i}(y_i)$ and scale $b_{|i}(y_i)$ functions are given by $a_{|i}(y_i) = \alpha_{|i} y_i$ and $b_{|i}(y_i) = y_i^{\beta_{|i}}$ where the vector
 159 constants $\alpha_{|i}$ and $\beta_{|i}$ take values of $\alpha_{j|i} \in [-1, 1]$ and $\beta_{j|i} \in (-\infty, 1)$, respectively, for all $j \neq i$. Finally,
 160 the dependence structure is described by the multivariate semi-parametric regression model:

$$161 \quad Y_{-i} = \alpha_{|i} y_i + y_i^{\beta_{|i}} Z_{|i} \text{ for } Y_i = y_i > u_{Y_i}, \quad i = 1, \dots, d.$$

162 The above equation expresses the behaviour of the vector variable Y , excluding the element of Y_i when
 163 it takes a large value. The dependence between the variables Y_i and Y_j is explained by the constant $\alpha_{j|i}$.
 164 Positive values indicate a positive relationship. The constant $\beta_{j|i}$ incorporates the changes in the
 165 variability of Y_j as Y_i increases. Details on estimating the dependence parameters are given in Heffernan
 166 and Tawn (2004) and Keef et al. (2013).

167 To obtain randomly generated samples of $X|X_i > u_{X_i}$, we adopted the following procedure. Initially,
 168 samples of Y_i from the Laplace distribution are simulated conditional on it exceeding its cumulative
 169 probability corresponding to $F_{X_i}(u_{X_i})$. Similarly, samples of random observations of $Z_{|i}$ are drawn
 170 from its estimated distribution $\hat{G}_{|i}$. Then, using the semi-parametric model, we obtain $Y_{-i} = \hat{\alpha}_{|i}y_i +$
 171 $y_i^{\hat{\beta}_{|i}}Z_{|i}$ and transform the vector $Y = (Y_{-i}, Y_i)$ to the originally distributed $X = (X_{-i}, X_i)$ by the inverse
 172 transformation.

173 2.4 Extreme Learning Machine (ELM)

174 The ELM is a data-driven method developed by Huang et al. (2006) that has been used effectively for
 175 streamflow forecasting (e.g., Deo and Şahin, 2016; Yaseen et al., 2016). Compared to other common
 176 ANN techniques, it has the advantages of fast learning speed and is characterised by improved
 177 performance in terms of commonly encountered problems, such as over-fitting and the effect of local
 178 minima. The model has a three-layer structure with one input, one hidden and a single output layer and
 179 can be expressed mathematically as:

$$180 \quad \sum_{i=1}^{\Lambda} B_i h_i(m_i \cdot x_t + n_i) = z_t$$

181 where Λ is the total number of nodes, B are the estimated weights between the nodes of the hidden and
 182 output layers, and $h(m, n, x)$ is the activation function with weights $m_i \in \mathfrak{R}^d$, biases $n_i \in \mathfrak{R}$ and the
 183 explanatory variable of the training dataset $x_t \in \mathfrak{R}^d$. Here, i and d denote the index of a specific hidden
 184 neuron (HN) and the number of input neurons, respectively, and Z is the model output.

185 Initially, the ELM model selects the input weights and hidden layer biases at random, and then
 186 calculates the output weights using a least squares method instead of adjusting them iteratively (see
 187 Chen et al. 2018 for details). Once the output weights \hat{B} have been estimated, forecasts are obtained by
 188 substituting the training dataset x_t with the testing one. The number of HNs in the hidden layer and the
 189 activation function are the only parameters that need to be pre-defined. The optimal number of HNs is
 190 a trade-off between generalization ability and network complexity. A highly complex model with too
 191 many HNs can lead to over-fitting, whereas a decreased number of HNs can result in a model that is
 192 too simple to capture non-linear relationships. The optimal number of HNs is problem-dependent and
 193 is frequently determined empirically (Huang et al., 2006; Sun et al., 2008). In this research, the number
 194 of HNs was increased iteratively from 1 to 100 and the network structure that provided the smallest
 195 RMSE of the training procedure was selected.

196 2.5 Application and Evaluation

197 A jackknife evaluation procedure (Miller, 1964; Shao and Tu, 1995) was applied to assess the
 198 performance of the proposed hybrid approach. It is a leave-one-out resampling technique without
 199 random replacement where one observation or a fixed subset of the dataset is omitted iteratively. The
 200 main strengths of the jackknife method are that model accuracy is independent of the calibration data
 201 and the loss in the sample data information is minimal (McCuen, 2005).

202 As stated previously, peak events are defined as flow above a certain threshold of the PBM simulated
 203 data. At each iteration, one peak flow event (measured and simulated) was left out of the dataset. This
 204 event constitutes the testing dataset and the rest of the data the training dataset, and the CEM and the

205 ELM were fitted to the latter. The dependence behavior of measured peaks conditional on the PBM
206 simulated, above a certain threshold, was configured by the CEM. From the fitted CEM, 50,000
207 stochastic simulations were obtained for both the observed X_j (pseudo-observations) and the PBM
208 simulated X_i variables (pseudo-PBM simulated). From the total set of random simulations of the
209 conditioning variable X_i , the ones with the smallest difference (≤ 0.1) from the maximum PBM
210 simulated peak of the testing sample, which was left out of the training dataset, were considered. As
211 CEM provides pairs of simulated data according to their dependence structure, the corresponding
212 random simulations of X_j (pseudo-observations) were then obtained. By calculating their median value,
213 a forecast of the maximum flow during an event was obtained and compared to the maximum measured
214 and PBM simulated peak excess of the testing dataset.

215 The ELM model was trained using PBM simulated data as inputs and measured data as outputs of the
216 training dataset. Based on the trained ELM model, flow forecasts were then obtained using the PBM
217 simulated flow of the testing sample as explanatory variable, except for the maximum. Consequently,
218 peaks smaller than the cluster maxima were forecasted by the ELM and the CEM was used only to
219 forecast maximum flows. The application of the ELM model alone on all the peaks was also performed
220 in experimentation and its performance compared to the CEM for the maximum flows. At the next
221 iteration, a different peak flow event was omitted from the training dataset for testing purposes and the
222 same process was repeated for all peaks.

223 This procedure was performed initially for peaks above the threshold that corresponds to the start of
224 the region of stability of shape and modified scale parameters. However, in order to investigate the
225 effect of threshold selection on the proposed methodology, the above-mentioned procedure was
226 repeated for different thresholds. The considered thresholds were set as a range from the minimum that
227 resulted from the application of threshold stability method, up to the 95th quantile of the PBM simulated
228 flow. Higher thresholds resulted in data scarcity that did not allow the models to be fitted satisfactorily.
229 All the above-mentioned steps are presented diagrammatically in Figure 1.

230 To assess the accuracy of the peak flow forecasts for each threshold, a set of indices was calculated.
231 More specifically, the mean absolute error (MAE), the normalized root mean square error (NRMSE),
232 the percentage BIAS (PBIAS), the Nash-Sutcliffe efficiency (NSE), the index of agreement (d) and the
233 Kling-Gupta Efficiency (KGE) were computed using the following equations:

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |\hat{z}_i - z_i|$$

$$\text{NRMSE} = 100 \frac{\sqrt{\frac{1}{N} \sum_{i=1}^N (\hat{z}_i - z_i)^2}}{z_{\max} - z_{\min}}$$

$$\text{PBIAS} = 100 \frac{\sum_{i=1}^N (\hat{z}_i - z_i)}{\sum_{i=1}^N z_i}$$

$$\text{NSE} = 1 - \frac{\sum_{i=1}^N (\hat{z}_i - z_i)^2}{\sum_{i=1}^N (z_i - \bar{z}_i)^2}$$

$$d = 1 - \frac{\sum_{i=1}^N (\hat{z}_i - z_i)^2}{\sum_{i=1}^N (|\hat{z}_i - \bar{z}_i| + |z_i - \bar{z}_i|)^2}$$

$$\text{KGE} = 1 - \sqrt{(r - 1)^2 + \left(\frac{\sigma_{\hat{z}}}{\sigma_z} - 1\right)^2 + \left(\frac{\bar{\hat{z}}}{\bar{z}} - 1\right)^2}$$

where \hat{z}_i are the simulated (or predicted) values, z_i are the measurements (or observed values), \bar{z}_i is the mean of the measured values, r is the Pearson product-moment correlation coefficient (between \hat{z}_i and z_i) and σ is the standard deviation. The optimal value of the error indices (MAE, NRMSE and PBIAS) is zero and the smaller are the values, the more accurate are the simulations. NSE (Nash and Sutcliffe, 1970) takes values from $-\infty$ to 1, where one corresponds to a perfect match between simulated and measured values, zero indicates that model simulations are as accurate as the mean of the measured values and a negative value indicates that the mean of the measured values is a more accurate predictor than the model. The index of agreement, d is defined in the range of zero to one, where again one represents the perfect model and zero no agreement at all. KGE incorporates r , the ratio between the means of the measurements and the simulations, and the variability ratio. KGE takes the same value range as NSE.

3 Study Site and Data

3.1 Study site

The flow discharge data used in this research were measured at the North Wyke Farm Platform (NWFP). The NWFP is a farm-scale experiment established in 2010 in the southwest of England (50°46'10"N, 3°54'05"W) to support research into sustainable grassland livestock systems (Orr et al., 2016). The platform comprises three independent small farms, each 21 ha in size. Each farm is divided into five sub-catchments, with some sub-catchments consisting of more than one field. The platform monitors routinely water run-off and water chemistry in each of the 15 sub-catchments, together with other primary data collections (e.g. greenhouse gas emissions) so that each farming system can be evaluated according to its level of sustainability (Takahashi et al., 2018). For the period 1985-2015, the average annual temperature at North Wyke ranges from 6.8 to 13.4 °C and the average annual rainfall is 1033 mm. The platform has an altitude range of 120–180 m above sea level. Soil texture

263 consists of a slightly stony clay loam topsoil (about 36% clay) above a mottled stony clay (about 60%
264 clay). The subsoil is impermeable to water and during rain events most of the excess water moves by
265 surface and sub-surface lateral flow towards the drainage system described below.

266 Each of the 15 sub-catchments (inset in Figure 2) are hydrologically isolated through a combination of
267 topography and a network of French drains (800-mm deep trenches) which ensure that the total runoff
268 is channelled to instrumented flumes, measuring water discharge and its chemistry with a 15 minute
269 temporal frequency since October 2012. The runoff from each sub-catchment is measured through a
270 combination of primary and secondary flow devices. The primary devices are H-type flumes
271 (TRACOM Inc., Georgia, USA) with capacity designed for a 1-in-50-year storm event (in respect of
272 data preceding 2010). The specific design of the H-type flume facilitates the accurate measurement of
273 both low and high flows and is relatively self-cleaning since it allows the ready passage of sediment
274 and particulate matter. A secondary flow measurement device (OTT hydromet, Loveland, CO., USA)
275 is used to measure the water height within the flume and convert it to discharge rate using flume-
276 specific formulas which depend on water height. The flow is generated only from rainfall as the fields
277 are not irrigated. Each sub-catchment also monitors precipitation and soil moisture every 15 minutes.

278 Platform data acquired from October 2011 to July 2013, represent a baseline period where all farm
279 fields were categorized as permanent pasture and received identical rates of inorganic fertilizers and
280 farmyard manure. From July 2013 to July 2015, two of the three farms entered a transition phase and
281 were ploughed and reseeded progressively with different types of pasture; specifically, a mixture of
282 white clover and high sugar perennial ryegrass, and sugar perennial ryegrass only. Thus, two farms
283 entered fully a post-baseline period in July 2015.

284 For this research, we used flow discharge (from April 2013 to February 2016) measured at sub-
285 catchment 6 of the permanent pasture farm (Figure 2), which consists of a single field (Golden Rove).
286 This field was chosen because, as part of the permanent pasture farm, it would not have been ploughed
287 and reseeded during the period of study (which would affect various processes, such as runoff).

288 **3.2 Choice of process-based model (PBM)**

289 For this research, we used the ‘SPACSYS’ model to simulate the flow discharge for sub-catchment 6
290 of the NWFP over the period of interest. The SPACSYS model is a process-based, field-scale model
291 which simulates key agricultural processes such as plant growth and development, soil Carbon and
292 Nitrogen (N) cycling, water dynamics and heat transformation (Wu et al., 2007) (see Figure 1). The
293 main processes concerning plant growth are assimilation, respiration, water and N uptake, partitioning
294 of photosynthate and N, N-fixation for legume plants and root growth. The Richards equation for water
295 potential is used in SPACSYS to simulate water redistribution in a soil profile. Site-specific input data
296 for the simulations include daily weather variables from the North Wyke site, soil properties, field and
297 grass management (e.g., fertiliser application dates and composition, reseeded, grazing and cutting
298 dates), and initialization of the state variables (standing biomass and root distribution, soil water and
299 temperature distribution). Previous simulations of water runoff, soil moisture and other agricultural
300 processes for sub-catchment 6 of the NWFP using SPACSYS can be found in Liu et al. (2018), where
301 a detailed explanation on the SPACSYS calibration is given.

302 **4 Results**

303 **4.1 Comparison of measured flow data with PBM simulations**

304 The plotted time-series of measured and PBM simulated flow (Figure 3), shows that the simulation
 305 appears to capture well the general behaviour of the process at low flows. However, it tends to under-
 306 predict the high flows and over-predict the medium ones. This is confirmed by the corresponding
 307 scatterplot (Figure 4) where many values in the range 5-10 mm d⁻¹ are below the 1-to-1 line and, thus,
 308 the simulated flow is greater than that measured. A non-linear locally weighted regression fit (i.e. a
 309 Loess smoother, see Cleveland, 1979), to the measured and simulated data is also given to help
 310 illustrate this behaviour.

311 4.2 Threshold selection

312 The shape and modified scale parameters estimated using the method of Curceac et al. (2020) indicated
 313 very similar threshold choices, in regions where the parameters remained relatively stable for
 314 increasing threshold candidates (Figure 5). The minimum threshold according to the shape parameter
 315 is 3.96 mm d⁻¹ and according to the modified scale parameter, 3.88 mm d⁻¹. These thresholds were
 316 estimated based on the PBM simulated flow (as described above), and the same thresholds were used
 317 for the observed peaks. Diagnostics, such as QQ plots of the empirical and modelled distributions (not
 318 presented), indicated that the GPD provides a good fit to the excesses and can model satisfactorily the
 319 peaks above the threshold of 3.88 mm d⁻¹, which was eventually selected. The range of thresholds
 320 above which the models were applied, was set from 3.88 mm d⁻¹ up to 6.41 mm d⁻¹, with the maximum
 321 corresponding to the 95th quantile of the PBM simulated flow.

322 4.3 Conditional Extreme Model (CEM) Fit

323 The diagnostics of the extreme dependence model (CEM) show a satisfactory fit (Figure 6). As stated
 324 in Section 2.3, one of the main assumptions of the model is that the residuals Z are independent of the
 325 conditioning variable (in this case, the PBM simulations). The pattern of both the initial and absolute
 326 values of the normalized residuals conforms approximately to a uniform distribution with no distinct
 327 pattern in the location or scatter of these residuals with the conditioning PBM simulations. The slight
 328 trend in the residuals Z for the lowest peaks of the conditioning variable might indicate that a higher
 329 threshold should be considered. The fitted quantiles of the conditional distribution of the dependent
 330 variable (measured data) conditional on the PBM simulated data (Figure 6, bottom) shows a good
 331 agreement between the data and the fitted quantiles, which capture the whole range of the scatter.
 332 Histograms of the scale and shape parameters (Figure 7) show that the measured and PBM simulated
 333 peaks have similar scale characteristics. However, the distribution of the measured peaks has a
 334 considerably heavier tail ($\xi_{obs} > \xi_{sim}$). The CEM simulated values of the dependent variable
 335 (measured data) along with the values of the conditional variable (PBM simulated data) (Figure 8)
 336 were obtained using the CEM with estimated dependence parameters of $\alpha = 0.44$ and $\beta = 0.59$. These
 337 parameters confirm that there is a positive dependence between the measured and the PBM simulated
 338 data, and that the measured data increase in variability as the values of the PBM simulations increase.

339 4.4 Hybrid model via CEM-ELM adjustments of PBM simulated data

340 To recap, this research applies the CEM for the maximum peaks, while the ELM model is used for the
 341 smaller peaks during a peak flow event as the ELM alone did not increase the accuracy of the maximum
 342 peaks (over that found with the PBM alone). For reference, error and agreement performance indices
 343 are given in Appendix A (Figure A1) for the three constituent models of the study hybrid (i.e. for PBM
 344 only, CEM only and ELM only), for predicting the maximum peaks.

345 The resultant hybrid simulations (or adjusted PBM simulations) for peak flow events above the
 346 minimum threshold of 3.88 mm d^{-1} are presented in Figure 9 together with the PBM simulated data
 347 and the measured data. The PBM most commonly under-predicts the largest peaks and over-predicts
 348 the ones preceding and following it. Use of the CEM captures the cluster maxima more accurately,
 349 which naturally depends on the value of the PBM simulation. In cases where the PBM over-predicts
 350 the maximum peak, the CEM leads to an even greater error. The ELM model addresses the fact that
 351 the PBM tends to over-predict the smaller peaks and, thus, provides hybrid forecasts of these peaks
 352 that are smaller and closer to the measured ones. The characteristics of the elements of the proposed
 353 methodology, in combination, results in improved characterization of the peak flow events, that tend
 354 to rise and fall more steeply (and realistically) than is found with the PBM simulations. Key exceptions
 355 arise for cases where the PBM over-predicts the whole event, as the hybrid compounds this over-
 356 prediction.

357 Error and agreement indices (Figure 10) provide an overall assessment of the proposed hybrid
 358 methodology for the same peak flow events (of Figure 9), but specifically just for instances of PBM
 359 simulations $> 3.88 \text{ mm d}^{-1}$. In general, the proposed hybrid approach is more accurate, as it results in
 360 smaller error indices and larger agreement indices than produced using the PBM alone, except for
 361 PBIAS, despite reductions in the other two error indices (MAE and NRMSE). Clearly, PBIAS is more
 362 reflective of how the hybrid can sometimes compound over-prediction. The greatest relative
 363 improvement was found in the KGE index, although both NSE and d also indicated improved
 364 agreement between observed and hybrid simulated values.

365 All of the results discussed above relate only to instances of PBM simulated flow values above the
 366 threshold of 3.88 mm d^{-1} , where the measured and hybrid simulated values directly correspond to. We
 367 compare now between *all* the measured water flow data, the PBM and hybrid simulations when above
 368 the selected threshold. The resultant plots of error (MAE and PBIAS only) and agreement (d and KGE
 369 only) indices against the magnitude of observed flow are given in Figure 11. The MAE is very small
 370 for both the PBM and the hybrid when comparing simulated flow with *all* the observed flow above the
 371 threshold. Increasing the observed flow threshold above which data are compared with the simulated
 372 data, results in a slower increase (with flow magnitude) in the MAE for the hybrid than for the PBM
 373 outputs. The hybrid approach also results in a significant decrease of the negative PBIAS with
 374 increasing peak flow, relative to the PBM. The agreement indices (d and KGE) similarly confirm this
 375 improvement found for the hybrid simulations over the PBM simulations.

376 All of the results discussed above refer to peak events above the threshold of 3.88 mm d^{-1} , as selected
 377 based on the GPD parameter stability plots (Figure 5). As a final step in the analysis, it is prudent to
 378 assess how threshold selection has an effect on the performance of the proposed methodology.
 379 Thresholds were set to range from 3.88 mm d^{-1} up to the 95th quantile of the PBM simulated flow (6.5
 380 mm d^{-1}). According to the calculated MAE indices, the hybrid model has a performance similar to the
 381 PBM when considering peak events above the threshold of 5.8 mm d^{-1} (Figure 12). This is not
 382 confirmed by the NRMSE which, however, shows a steep increase for the same threshold. PBIAS
 383 shows an overall increasing trend with some fluctuations in between. The agreement indices (Figure
 384 12) seem to be less sensitive to the threshold, although NSE shows an abrupt decrease when flow is
 385 higher than 5.8 mm d^{-1} . All the indices have the common characteristic of the consistent trend
 386 (increasing for error, decreasing for agreement) as the threshold increases, which could be attributed
 387 to the smaller samples of the data used for testing, in which the highest flow values dominate.

388 5 Discussion

389 The main motivation for developing the proposed hybrid approach was to forecast more accurately the
390 peak flows that are typically under-predicted using PBMs due to model over-generalisation or
391 smoothing. The analysis in this research was based on simulations obtained from the SPACSYS model.
392 SPACSYS has characteristics that can be considered as representative of the vast majority of PBMs
393 used for flow simulations and the hybrid approach presented is entirely general. However, the PBM
394 also exhibited other problems, such as over-predicting small and moderate flow values. This second
395 problem arises because the model (as for most PBMs) is calibrated implicitly to the *mean* of the
396 observed distribution through the careful choice and selection of model parameters. It should be noted,
397 however, that SPACSYS is not fitted or re-calibrated explicitly to external data.

398 Topological characteristics, such as the integrating effect of the catchment, could also contribute to this
399 behaviour. For example, large local slopes (that SPACSYS cannot represent) result in faster running
400 water which, combined with intense rainfall, may result in higher peak flows that are not captured by
401 SPACSYS. Over-predicted events are likely due to inaccurate representation of soil moisture,
402 topography and other soil properties at the within-field scale, since SPACSYS simulates at the field
403 scale (Liu et al., 2018). Despite these issues and the fact that our proposed hybrid approach was aimed
404 at under-predicted extreme flow events, the hybrid approach resulted in more accurate forecasts and
405 an increase in accuracy overall.

406 The CEM is usually used to describe the extreme dependence structure of the same variable at different
407 sites or of different variables at the same site. In this study, we used the CEM in a bivariate context to
408 model and link the same underlying state variable captured by different representational processes (i.e.,
409 direct measurement and PBM simulation of flow). The pseudo-observations obtained from the fitted
410 model and based on the conditioning variable were aggregated to a single value which was then
411 compared to the equivalent measured value. The same conditional simulations can be used to create
412 confidence intervals that correspond to various scenarios and allow flexibility in choosing values
413 according to the intended purpose.

414 In general, none of the applied criteria for the evaluation of the proposed hybrid method is sufficient
415 singly; each of the model performance indices have strengths and weaknesses. The agreement indices
416 are used mainly to investigate how accurately the model captures the dynamic of the temporal process.
417 The error indices capture differences between the total flow or the volume of the hydrograph.
418 Therefore, using both measures provides a more holistic evaluation of model performance. Since our
419 main objective was to evaluate the performance of the proposed hybrid method in predicting extreme
420 flows, the choice of the agreement indices is appropriate as they have been shown to be sensitive to
421 peaks (Krause et al., 2005).

422 Despite the promising results obtained from the proposed methodology, it has the limitation of being
423 tested for a specific case study site and for one PBM. Future research should, therefore, consider testing
424 this approach for other catchment sites with different characteristics, as data-driven models need to be
425 tested using a range of (large) datasets before applied in practice (Boulesteix et al., 2018;
426 Papacharalampous et al., 2019; Tyrallis et al., 2019). It would also be interesting to investigate whether
427 and how the performance of SPACSYS, and by extension, the proposed techniques, would be affected
428 by using forecasted weather variables as inputs instead of measured data to obtain the simulations. In
429 real case scenarios, the threshold is defined commonly based on pre-existing information. Due to the
430 nature of the NWFP experiment, it was not possible to define a threshold with physical meaning (e.g.
431 likely flooding) with which to evaluate the estimated threshold. The threshold defines the peak flow
432 events and consequently the training and testing datasets used in this research. Thus, it was not possible

433 to define a threshold based strictly on the training dataset only as would normally be the case. However,
434 we expect this to have a minimal effect on the results and not change the main conclusions drawn.

435 **Conclusions**

436 In this research, we used a data-driven machine learning model (ELM) and a semi-parametric
437 conditional model that stems from extreme value theory (CEM) to increase the accuracy of peak water
438 flow events simulated by a process-based model (PBM). The PBM most frequently under-predicted
439 the maximum flows during a peak event, for which the CEM was applied, and over-predicted flows
440 preceding and following it, for which the ELM was applied. The combined characteristics of the
441 proposed methodology in general resulted in more accurate forecasts and improved representation of
442 these peak events, according to several error and agreement indices. The detailed analysis undertaken
443 in this research was developed based on simulated flow data obtained from only one PBM and for
444 observed data at only one case study site. However, because of the general characteristics of the chosen
445 PBM and of the proposed hybrid methodology, it is anticipated that the proposed approach will be
446 suitable for a wide range of PBMs and water monitoring station schemes.

447 **Conflict of Interest**

448 The authors declare that the research was conducted in the absence of any commercial or financial
449 relationships that could be construed as a potential conflict of interest.

450 **Acknowledgements**

451 Rothamsted Research receives grant aided support from the Biotechnology and Biological Sciences
452 Research Council (BBSRC) of the United Kingdom. This research was funded by Rothamsted
453 Research and Lancaster Environment Centre, the BBSRC Institute Strategic Programme (ISP) grant,
454 “Soils to Nutrition” (S2N) grant numbers BBS/E/C/000I0320, BBS/E/C/000I0330 and the BBSRC
455 National Capability grant for the North Wyke Farm Platform grant number BBS/E/C/000J0100. The
456 authors wish to thank the Editor and two anonymous reviewers for their useful comments, which led
457 to considerable improvements to the paper.

458 **Data and Software Availability Statement**

459 All North Wyke Farm Platform datasets (<https://www.rothamsted.ac.uk/north-wyke-farm-platform>)
460 and the SPACSYS model (<https://www.rothamsted.ac.uk/rothamsted-spacsys-model>) are freely
461 available. R software (R Core Team, 2019) was used for the implementation of the statistical models.
462 The CEM was applied by using the texmex R package (Southworth et al., 2018), the elmNMRcpp R
463 package was used for the ELM model (Mouselimis and Gosso, 2018) and the indices were calculated
464 by using functions in the hydroGOF R package (Zambrano-Bigiarini, 2017).

465

466 **References**

- 467 Bates, B. C., Kundzewicz, Z. W., Wu, S. and Palutikof, J. P. (2008). Climate Change and Water.
468 Technical Paper of the Intergovernmental Panel on Climate Change, IPCC Secretariat, Geneva,
469 210 pp.
- 470 Bogner, K., Liechti, K. and Zappa, M. (2016). Post-Processing of Stream Flows in Switzerland with
471 an Emphasis on Low Flows and Floods, *Water*, 8 (4), 115. doi: 10.3390/w8040115.
- 472 Bogner, K., Liechti, K. and Zappa, M. (2017). Technical Note: Combining Quantile Forecasts and
473 Predictive Distributions of Streamflows, *Hydrology and Earth System Sciences*, 21 (11), 5493-
474 5502. doi: 10.5194/hess-21-5493-2017.
- 475 Boulesteix, A. L., Binder, H., Abrahamowicz, M. and Sauerbrei, W. (2018). On the Necessity and
476 Design of Studies Comparing Statistical Methods, *Biometrical Journal*, 60 (1), 216-218. doi:
477 10.1002/bimj.201700129.
- 478 Bouraoui, F., Grizzetti, B., Granlund, K., Rekolainen, S. and Bidoglio, G. (2004). Impact of Climate
479 Change on the Water Cycle and Nutrient Losses in a Finnish Catchment, *Climatic Change*,
480 66(1–2), 109-126. doi.org: 10.1023/B:CLIM.0000043147.09365.e3.
- 481 Box, G. E. P. and Jenkins, G. M. (1976). *Time Series Analysis: Forecasting and Control*, Revised
482 Edition. Holden-Day, San Francisco, CA.
- 483 Bradley, A. A., Habib, M. and Schwartz, S. S. (2015). Climate index weighting of ensemble streamflow
484 forecasts using a simple Bayesian approach, *Water Resources Research*, 51, 7382–7400. doi:
485 10.1002/2014WR016811.
- 486 Chen, L., Sun, N., Zhou, C., Zhou, J., Zhou, Y., Zhang, J. and Zhou, Q. (2018). Flood Forecasting
487 Based on an Improved Extreme Learning Machine Model Combined with the Backtracking
488 Search Optimization Algorithm, *Water*, 10(10), 1362. doi: 10.3390/w10101362.
- 489 Cleveland, W. S. (1979). Robust Locally Weighted Regression and Smoothing Scatterplots, *Journal of*
490 *the American Statistical Association*, 74(368), 829-836. doi: 10.2307/2286407.
- 491 Cloke, H. L. and Pappenberger, F. (2009). Ensemble Flood Forecasting: A Review, *Journal of*
492 *Hydrology*, 375(3), 613-626. doi: 10.1016/j.jhydrol.2009.06.005.
- 493 Collet, L., Beevers, L. and Prudhomme, C. (2017). Assessing the Impact of Climate Change and
494 Extreme Value Uncertainty to Extreme Flows across Great Britain, *Water*, 9(2), 103. doi:
495 10.3390/w9020103.
- 496 Curceac, S., Atkinson, P. M., Milne, A., Wu, L. and Harris, P. (2020). An Evaluation of Automated
497 GPD Threshold Selection Methods for Hydrological Extremes across Different Scales, *Journal*
498 *of Hydrology*, 585, 124845. doi: 10.1016/j.jhydrol.2020.124845.
- 499 Deo, R. C. and Şahin, M. (2016). An Extreme Learning Machine Model for the Simulation of Monthly
500 Mean Streamflow Water Level in Eastern Queensland, *Environmental Monitoring and*
501 *Assessment*, 188, 90. doi: 10.1007/s10661-016-5094-9.
- 502 Dogulu, N., López López, P., Solomatine, D. P., Weerts, A. H. and Shrestha, D. L. (2015). Estimation
503 of Predictive Hydrologic Uncertainty Using the Quantile Regression and UNEEC Methods and
504 Their Comparison on Contrasting Catchments, *Hydrology and Earth System Sciences*, 19 (7),
505 3181-3201. doi: 10.5194/hess-19-3181-2015.
- 506 Drees, H. and Janßen, A. (2017). Conditional Extreme Value Models: Fallacies and Pitfalls, *Extremes*,
507 20(4), 777–805. doi: 10.1007/s10687-017-0293-5.

- 508 Fathian, F., Mehdizadeh, S., Kozekalani A. S. and Safari, M. J. S. (2019). Hybrid Models to Improve
 509 the Monthly River Flow Prediction: Integrating Artificial Intelligence and Non-Linear Time
 510 Series Models, *Journal of Hydrology*, 575, 1200–1213. doi: 10.1016/j.jhydrol.2019.06.025.
- 511 Field, C. B., Barros, V., Stocker, T. F. and Dahe, Q. (2012). Managing the Risks of Extreme Events
 512 and Disasters to Advance Climate Change Adaptation: Special Report of the Intergovernmental
 513 Panel on Climate Change, Cambridge, Cambridge University Press. doi:
 514 10.1017/CBO9781139177245.
- 515 Heffernan, J. E. and Tawn, J. A. (2004). A Conditional Approach for Multivariate Extreme Values
 516 (with Discussion). *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*,
 517 66(3), 497-546. doi.org: 10.1111/j.1467-9868.2004.02050.x.
- 518 Huang, G. B., Zhu, Q. Y. and Siew, C. K. (2006). Extreme Learning Machine: Theory and
 519 Applications, *Neurocomputing, Neural Networks*, 70(1), 489-501. doi:
 520 10.1016/j.neucom.2005.12.126.
- 521 Keef, C., Papastathopoulos, I. and Tawn, J. A. (2013). Estimation of the Conditional Distribution of a
 522 Multivariate Variable given That One of Its Components Is Large: Additional Constraints for
 523 the Heffernan and Tawn Model, *Journal of Multivariate Analysis*, 115, 396-404. doi:
 524 10.1016/j.jmva.2012.10.012.
- 525 Kisi, O. and Cimen, M. (2011). A Wavelet-Support Vector Machine Conjunction Model for Monthly
 526 Streamflow Forecasting, *Journal of Hydrology*, 399(1), 132-140. doi:
 527 10.1016/j.jhydrol.2010.12.041.
- 528 Krause, P., Boyle, D. P. and Bäse, F. (2005). Comparison of Different Efficiency Criteria for
 529 Hydrological Model Assessment, *Advances in Geosciences*, 5, 89-97. doi: 10.5194/adgeo-5-
 530 89-2005.
- 531 Kundzewicz, Z. W., Mata, L. J., Arnell, N. W., Doll, P., Kabat, P., Jimenez, B. et al. (2007). Freshwater
 532 Resources and Their Management. In *Climate Change 2007: Impacts, Adaptation and*
 533 *Vulnerability. Contribution of Working Group II to the Fourth Assessment Report of the*
 534 *Intergovernmental Panel on Climate Change*, edited by M. L. Parry, O. F. Canziani, J. P.
 535 Palutikof, P. J. van der Linden, and C. E. Hanson, 173–210. Cambridge University Press.
- 536 Lamb, R., Keef, C., Tawn, J., Laeger, S., Meadowcroft, I., Surendran, S., Dunning, P. and Batstone,
 537 C. (2010). A New Method to Assess the Risk of Local and Widespread Flooding on Rivers and
 538 Coasts, *Journal of Flood Risk Management*, 3(4), 323-336. doi: 10.1111/j.1753-
 539 318X.2010.01081.x.
- 540 Lane, R. A., Coxon, G., Freer, J. E., Wagener, T., Johnes, P. J., Bloomfield, J. P., Greene, S., Macleod,
 541 C. J. A. and Reaney, S. M. (2019). Benchmarking the Predictive Capability of Hydrological
 542 Models for River Flow and Flood Peak Predictions across over 1000 Catchments in Great
 543 Britain, *Hydrology and Earth System Sciences*, 23(10), 4011-4032. doi: 10.5194/hess-23-4011-
 544 2019.
- 545 Li, W., Duan, Q., Miao, C., Ye, A., Gong, W. and Di, Z. (2017). A Review on Statistical Postprocessing
 546 Methods for Hydrometeorological Ensemble Forecasting. *Wiley Interdisciplinary Reviews*,
 547 *Water*, 4(6): e1246. doi: 10.1002/wat2.1246.
- 548 Li, X. Q., Chen, J., Xu, C. Y., Li, L. and Chen, H. (2019). Performance of Post-Processed Methods in
 549 Hydrological Predictions Evaluated by Deterministic and Probabilistic Criteria, *Water*
 550 *Resources Management*, 33(9), 3289-3302. doi: 10.1007/s11269-019-02302-y.

- 551 Lima, A. R., Cannon, A. J. and Hsieh, W. W. (2015). Nonlinear Regression in Environmental Sciences
 552 Using Extreme Learning Machines: A Comparative Evaluation, *Environmental Modelling &*
 553 *Software*, 73, 175-188. doi: 10.1016/j.envsoft.2015.08.002.
- 554 Liu, Y., Li, Y., Harris, P., Cardenas, L. M., Dunn, R. M., Sint, H., Murray, P. J., Lee, M. R. F. and Wu,
 555 L. (2018). Modelling Field Scale Spatial Variation in Water Run-off, Soil Moisture, N₂O
 556 Emissions and Herbage Biomass of a Grazed Pasture Using the SPACSYS Model, *Geoderma*,
 557 315, 49-58. doi: 10.1016/j.geoderma.2017.11.029.
- 558 López López, P., Verkade, J. S., Weerts, A. H. and Solomatine, D. P. (2014). Alternative
 559 Configurations of Quantile Regression for Estimating Predictive Uncertainty in Water Level
 560 Forecasts for the Upper Severn River: A Comparison. *Hydrology and Earth System Sciences*,
 561 18(9), 3411-3428. doi: 10.5194/hess-18-3411-2014.
- 562 McCuen R. H. (2005). Accuracy Assessment of Peak Discharge Models, *Journal of Hydrologic*
 563 *Engineering*, 10, (1), 16-22. doi: 10.1061/(ASCE)1084-0699(2005)10:1(16).
- 564 Mendes, B. V. de M. and Pericchi, L. R. (2009). Assessing Conditional Extremal Risk of Flooding in
 565 Puerto Rico, *Stochastic Environmental Research and Risk Assessment*, 23(3), 399-410. doi:
 566 10.1007/s00477-008-0220-z.
- 567 Miller, R. G. (1964). A Trustworthy Jackknife, *The Annals of Mathematical Statistics*, 35(4), 1594-
 568 1605. doi: 10.1214/aoms/1177700384.
- 569 Mouselimis, L. and Gosso, A. (2018). elmNNRcpp: The Extreme Learning Machine Algorithm. R
 570 package version 1.0.1. <https://CRAN.R-project.org/package=elmNNRcpp>
- 571 Nash, J. E. and Sutcliffe, J. V. (1970). River Flow Forecasting through Conceptual Models Part I - A
 572 Discussion of Principles, *Journal of Hydrology*, 10, (3): 282-290. doi: 10.1016/0022-
 573 1694(70)90255-6.
- 574 Orr, R. J., Murray, P. J., Eyles, C. J., Blackwell, M. S. A., Cardenas, L. M., Collins, A. L. et al. (2016).
 575 The North Wyke Farm Platform: effect of temperate grassland farming systems on soil moisture
 576 contents, runoff and associated water quality dynamics, *European Journal of Soil Science*, 67,
 577 374–385. doi: 10.1111/ejss.12350.
- 578 Papacharalampous, G., Tyralis, H., Langousis, A., Jayawardena, A. W., Sivakumar, B., Mamassis, N.,
 579 Montanari, A. and Koutsoyiannis, D. (2019). Probabilistic Hydrological Post-Processing at
 580 Scale: Why and How to Apply Machine-Learning Quantile Regression Algorithms, *Water*,
 581 11(10), 2126. doi: 10.3390/w11102126.
- 582 Quilty, J., Adamowski, J. and Boucher, M. A. (2019). A Stochastic Data-Driven Ensemble Forecasting
 583 Framework for Water Resources: A Case Study Using Ensemble Members Derived From a
 584 Database of Deterministic Wavelet-Based Models, *Water Resources Research*, 55(1), 175-202.
 585 doi: 10.1029/2018WR023205.
- 586 Raftery, A. E., Gneiting, T., Balabdaoui, F. and Polakowski, M. (2005). Using Bayesian Model
 587 Averaging to Calibrate Forecast Ensembles, *Monthly Weather Review*, 133(5), 1155-1574. doi:
 588 10.1175/MWR2906.1.
- 589 Roulin, E. and Vannitsem, S. (2011). Postprocessing of Ensemble Precipitation Predictions with
 590 Extended Logistic Regression Based on Hindcasts, *Monthly Weather Review*, 140(3), 874-888.
 591 doi: 10.1175/MWR-D-11-00062.1.
- 592 Scarrott, C. and MacDonald, A. (2012). A Review of Extreme Value Threshold Es-Timation and
 593 Uncertainty Quantification, *REVSTAT–Statistical Journal*, 10(1), 33-60.

- 594 Shao, J. and Tu, D. (1995). *The Jackknife and Bootstrap*, Springer Series in Statistics, New York,
595 Springer-Verlag. doi: 10.1007/978-1-4612-0795-5.
- 596 Sikorska, A. E., Montanari, A. and Koutsoyiannis, D. (2015). Estimating the Uncertainty of
597 Hydrological Predictions through Data-Driven Resampling Techniques, *Journal of Hydrologic*
598 *Engineering*, 20(1), A4014009. doi: 10.1061/(ASCE)HE.1943-5584.0000926.
- 599 Southworth, H., Heffernan J. E. and Metcalfe, P. D. (2018). *texmex: Statistical modelling of extreme*
600 *values*. R package version 2.4.2.
- 601 Sun, Z. L., Choi, T. M., Au, K. F. and Yu, Y. (2008). Sales Forecasting Using Extreme Learning
602 Machine with Applications in Fashion Retailing, *Decision Support Systems*, 46(1), 411-419.
603 doi: 10.1016/j.dss.2008.07.009.
- 604 Takahashi, T., Harris, P. M., Blackwell, S. A., Cardenas, L. M., Collins, A. L., Dungait, J. A. J.,
605 Hawkins, J. M. B. et al. (2018). Roles of Instrumented Farm-Scale Trials in Trade-off
606 Assessments of Pasture-Based Ruminant Production Systems, *Animal*, 12(8),1766-1776. doi:
607 10.1017/S1751731118000502.
- 608 Thibault, K. M. and Brown, J. H. (2008). Impact of an Extreme Climatic Event on Community
609 Assembly, *Proceedings of the National Academy of Sciences*, 105(9), 3410-3415. doi:
610 10.1073/pnas.0712282105.
- 611 Toth, E., Montanari, A. and Brath, A. (1999). Real-Time Flood Forecasting via Combined Use of
612 Conceptual and Stochastic Models, *Physics and Chemistry of the Earth, Part B: Hydrology,*
613 *Oceans and Atmosphere*, 24(7), 793-798. doi: 10.1016/S1464-1909(99)00082-9.
- 614 Tyralis, H., Papacharalampous, G., Burnetas, A. and Langousis, A. (2019). Hydrological Post-
615 Processing Using Stacked Generalization of Quantile Regression Algorithms: Large-Scale
616 Application over CONUS, *Journal of Hydrology*, 577, 123957. doi:
617 10.1016/j.jhydrol.2019.123957.
- 618 Wijayarathne, D. B. and Coulibaly, P. (2020). Identification of Hydrological Models for Operational
619 Flood Forecasting in St. John's, Newfoundland, Canada, *Journal of Hydrology: Regional*
620 *Studies*, 27, 100646. doi: 10.1016/j.ejrh.2019.100646.
- 621 Wu, L., McGechan, M. B. McRoberts, N., Baddeley, J. A. and Watson, C. A. (2007). SPACSYS:
622 Integration of a 3D Root Architecture Component to Carbon, Nitrogen and Water Cycling-
623 Model Description, *Ecological Modelling*, 200(3), 343-359. doi:
624 10.1016/j.ecolmodel.2006.08.010.
- 625 Yaseen, Z. M. Jaafar, O., Deo, R. C., Kisi, O., Adamowski, J., Quilty, J. and El-Shafie, A. (2016).
626 Stream-Flow Forecasting Using Extreme Learning Machines: A Case Study in a Semi-Arid
627 Region in Iraq, *Journal of Hydrology*, 542, 603-614. doi: 10.1016/j.jhydrol.2016.09.035.
- 628 Yaseen, Z. M., Sulaiman, S. O., Deo, R. C. and Chau, K. W. (2019). An Enhanced Extreme Learning
629 Machine Model for River Flow Forecasting: State-of-the-Art, *Practical Applications in Water*
630 *Resource Engineering Area and Future Research Direction*, *Journal of Hydrology*, 569, 387-
631 408. doi: 10.1016/j.jhydrol.2018.11.069.
- 632 Zambrano-Bigiarini, M. (2017). *hydroGOF: Goodness-of-fit functions for comparison of simulated*
633 *and observed hydrological time series*R package version 0.3-
634 [10.http://hzambran.github.io/hydroGOF/](http://hzambran.github.io/hydroGOF/). DOI:10.5281/zenodo.840087.

- 635 Zheng, F., Westra, S., Leonard, M. and Sisson, S. A. (2014). Modeling Dependence between Extreme
636 Rainfall and Storm Surge to Estimate Coastal Flooding Risk, *Water Resources Research*, 50(3),
637 2050-2071. doi: 10.1002/2013WR014616.
- 638 Zhou, J., Peng, T., Zhang, C. and Sun, N. (2018). Data Pre-Analysis and Ensemble of Various Artificial
639 Neural Networks for Monthly Streamflow Forecasting, *Water*, 10(5), 628. doi:
640 10.3390/w10050628.
- 641