

# The Promise of Diversity: Distribution-based Hydrologic Model Evaluation and Diagnostics

Jasper A. Vrugt\*† 

March 26, 2023

---

\*Department of Civil and Environmental Engineering, University of California, Irvine, California, USA.  
†Corresponding author: [jasper@uci.edu](mailto:jasper@uci.edu)

## Highlights

1. This paper advocates the use of Monte Carlo-based simulation distributions for hydrologic model evaluation and model diagnostics.
2. Distribution forecasts coalesce information about model behavior, robustness, sensitivity and uncertainty that is not available in single-valued model output.
3. *Strictly proper* scoring rules provide a rigorous information-theoretic underpinning to hydrologic model evaluation and diagnostics.
4. The Bregman divergence ensures *strict propriety* of scoring rules and their associated divergence functions.
5. Propriety and elicability offer useful working paradigms for the development, use and standardization of hydrologic scoring functions and scoring rules.
6. We present divergence scores for flood frequency analysis and flow duration and recession curves.

# <sup>1</sup> Abstract

<sup>2</sup> Distribution forecasts over future quantities or events are routinely made in hydrology but evade model  
<sup>3</sup> evaluation by a need for ensemble simulation and lack of knowledge on how to properly evaluate simulation  
<sup>4</sup> distributions against data. Predictive distributions  $P$  derived from (quasi)-Bayesian methods are usually  
<sup>5</sup> traded for a (likelihood-weighted) mean or median prediction to accommodate error measures (scoring  
<sup>6</sup> functions) such as the mean absolute error or the mean squared error that depend both on the point  
<sup>7</sup> forecasts and the realizing observations. Point in case is the so-called KG efficiency or KGE of *Gupta et al.*  
<sup>8</sup> (2009) and improvements thereof (*Lamontagne et al.*, 2020), which have rapidly gained popularity as  
<sup>9</sup> scoring functions among hydrologists as alternative to the infamous *Nash and Sutcliffe* (1970) efficiency,  
<sup>10</sup> but are equally exclusive in how they quantify model performance using only single-valued simulated  
<sup>11</sup> output of the quantities of interest. This paper advocates the use of simulation distributions for hydrologic  
<sup>12</sup> model evaluation and model diagnostics. Distribution evaluation is supported by information-theoretic  
<sup>13</sup> arguments and puts into modeling practice the social justice narrative of diversity, equity and inclusion  
<sup>14</sup> for different simulations. We discuss past developments that led to the current state-of-the-art of  
<sup>15</sup> forecast verification in hydrology and bring to the fore scoring rules for model evaluation and diagnostics.  
<sup>16</sup> *Strictly proper* scoring rules condense a distribution forecast to a single reward value for the materialized  
<sup>17</sup> outcome(s) and have a strong underpinning in statistical, decision and information theory. We review  
<sup>18</sup> scoring rules for dichotomous and categorical events, quantiles (intervals) and density forecasts, discuss  
<sup>19</sup> the importance of scoring rule propriety and address diagnostic aspects of a distribution forecast such as  
<sup>20</sup> sharpness, reliability and entropy. The usefulness and power of scoring rules is demonstrated on simple  
<sup>21</sup> benchmark problems and discharge distributions simulated with conceptual watershed models using  
<sup>22</sup> Generalized Likelihood Uncertainty Estimation and Bayesian model averaging. We also link scoring rules  
<sup>23</sup> to model diagnostics and present *strictly proper* divergence scores for flood frequency analysis and flow  
<sup>24</sup> duration and recession curves. Scoring rules offer a rigorous information-theoretic underpinning to model  
<sup>25</sup> evaluation and diagnostics and provide statistically principled means for (Bayesian) model selection and  
<sup>26</sup> the analysis of hydrograph functionals, flood frequencies and extreme events.

<sup>27</sup> **Keywords:** Ensemble prediction, distribution forecast, information theory, scoring rules, divergence score,  
<sup>28</sup> propriety, sharpness, reliability, uncertainty, entropy, integral transform, logarithmic score, continuous  
<sup>29</sup> ranked probability score, recession analysis, flow duration curve, signatures, watershed models

# <sup>30</sup> 1 Introduction and Scope

<sup>31</sup> The topic of model evaluation has received considerable attention in the hydrologic and water resources  
<sup>32</sup> literature over the past decades. Model evaluation is an integral part of the model development process  
<sup>33</sup> and involves comparing simulated system behaviour with observations in pursuit of a qualitative and/or  
<sup>34</sup> quantitative understanding of their similarities and differences and how well the model approximates  
<sup>35</sup> reality according to some error measure. This process must acknowledge differences in extent, support  
<sup>36</sup> and spacing of modeled and observed quantities (*Grayson and Blöschl*, 2001). We follow *Gneiting* (2011)  
<sup>37</sup> and use the terminology of a scoring function on a general sample space  $\Omega$  for an error measure.

<sup>38</sup> **Definition 1.** *A scoring function is any real-valued function  $s : \Omega \times \Omega \rightarrow \mathbb{R}$  where  $s(y, \omega)$  represents the  
loss or penalty when the point forecast  $y \in \Omega$  is issued and the observation  $\omega \in \Omega$  materializes.*

<sup>40</sup> Thus, scoring functions such as the ubiquitous squared error,  $s_{SE}(y, \omega) = (\omega - y)^2$ , absolute error  
<sup>41</sup>  $s_{AE}(y, \omega) = |\omega - y|$ , and *Nash and Sutcliffe* (1970) efficiency,  $s_{NS}(y, \omega) = 1 - (\omega - y)^2 / (\omega - \mu_\omega)^2$ , measure  
<sup>42</sup> the performance of a point forecast  $y$ , where  $\mu_\omega$  is the mean of the verifying data and  $\epsilon = \omega - y$  is  
<sup>43</sup> the so-called residual. The  $s_{SE}(y, \omega)$  and  $s_{AE}(y, \omega)$  scoring functions easily generalize to  $n$ -vectors of  
<sup>44</sup> forecasts and observations. We use the error terminology for  $s_{SE}(y, \omega)$  and  $s_{AE}(y, \omega)$  though it would be  
<sup>45</sup> more precise to use the wording of squared and absolute residual (*Vrugt and de Oliveira*, 2022). For the  
<sup>46</sup> purpose of this discussion, we shall classify the research on model evaluation into two different groups  
<sup>47</sup> including (i) residual-based approaches and (ii) non-residual-based methods. Methods in this second  
<sup>48</sup> group may still use residuals to quantify model performance but these are usually not discharge residuals  
<sup>49</sup> but residuals of some hydrograph functional.

<sup>50</sup> The residual-based approach has its roots in linear regression theory and determines model performance  
<sup>51</sup> using goodness-of-fit measures of simulated and measured watershed behavior. This includes the use  
<sup>52</sup> of (a) formal objective functions, likelihood functions and summary metrics that result from rigorous  
<sup>53</sup> application of statistical principles to the assumed probabilistic properties of the residuals within the  
<sup>54</sup> context of weighted and/or generalized least squares (*Tasker*, 1980; *Stedinger and Tasker*, 1985; *Kavetski*,  
<sup>55</sup> *et al.*, 2006a,b), maximum likelihood and Bayesian estimation (*Sorooshian and Dracup*, 1980; *Kuczera*,  
<sup>56</sup> 1983; *Bates and Campbell*, 2001; *Schoups and Vrugt*, 2010; *Scharnagl et al.*, 2015; *Ammann et al.*, 2019;  
<sup>57</sup> *Vrugt et al.*, 2022) approximate Bayesian computation (*Nott et al.*, 2012; *Vrugt and Sadegh*, 2013a; *Sadegh*

58 and Vrugt, 2013) and information-theoretic approaches (Neuman, 2003; Ye *et al.*, 2008; Weijs *et al.*,  
59 2010a; Lu *et al.*, 2011; Schöniger *et al.*, 2014; Pachepsky *et al.*, 2016; Volpi *et al.*, 2017), (b) informal  
60 metrics of the quality-of-fit and pseudo-likelihood functions within the context of model training (Nash  
61 and Sutcliffe, 1970; Gupta *et al.*, 2009; Pool *et al.*, 2018; Knoben *et al.*, 2019; Lamontagne *et al.*, 2020;  
62 Schwemmle *et al.*, 2021), informal Bayesian approaches (Beven and Binley, 1992; Freer *et al.*, 1996;  
63 Beven and Freer, 2001) and multi-criteria model calibration (Gupta *et al.*, 1998; Boyle *et al.*, 2000), and  
64 (c) tolerable ranges of the residuals within the context of limits of acceptability (Beven, 2006; Vrugt and  
65 Beven, 2018), regional sensitivity analysis (Spear and Hornberger, 1980; Spear *et al.*, 2020), dynamic  
66 identifiability analysis (Wagener *et al.*, 2003) and the parameter identification method based on the  
67 localization of information (Vrugt *et al.*, 2002). Let it be evident that residual-based methods are often  
68 used for parameter estimation purposes so as to train/calibrate models to mimic the observed data  
69 as closely and consistently as possible. While it is desirable that the model can reproduce historical  
70 observations, the aggregation of the residuals into a single performance index implies a significant loss  
71 of information about specific system behaviours. As a result, residual-based methods provide only  
72 limited guidance on model improvement. If the model is expected to match a certain functional of the  
73 hydrograph, it is critical that the scoring function be consistent for it, in the sense that the expected score  
74 is maximized (or minimized, if appropriate) when following the directive. Functionals that incentivize a  
75 truthful model description are called elicitable in decision-theory.

76 Non-residual-based methods are usually much more geared towards hypothesis testing and learning  
77 by exploiting hydrological context and theory, though they may be used for model training as well.  
78 These methods do not work with the residuals of simulated and measured watershed behavior but rather  
79 quantify model performance using multiple different hydrologic functionals of the streamflow hydrograph.  
80 Examples of such numeral descriptors of the catchment response to rainfall are the runoff ratio, baseflow  
81 index and flow duration curve. Shamir *et al.* (2005) laid the foundation of this approach in their work  
82 on hydrograph indices - and this approach advanced and matured further into what is now known as  
83 model diagnostics (Gupta *et al.*, 2008). As hydrologic functionals of the streamflow hydrograph are  
84 elicitable, discrepancies between measured and simulated functionals (signatures) are symptoms of model  
85 malfunctioning, and if able to relate each functional to a specific process, will guide model improvement in  
86 a more meaningful and systematic way. At its strongest, diagnostic evaluation will point clearly towards

87 the aspects of the model that need improvement, and give guidance towards the manner of improvement  
88 (*Gupta et al., 2008; Yilmaz et al., 2008; Westerberg et al., 2011*).

89 Our classification of model evaluation methods into residual and non-residual based approaches suffices  
90 for the purpose of this paper but is not complete. One may discern a group of residual-based model  
91 evaluation strategies with a diagnostic intent. Examples of this third group are model evaluation methods  
92 that couple ubiquitous scoring functions with wavelets (*Rathinasamy et al., 2014*), self-organizing maps  
93 (*Reusser et al., 2009*), interval deviation (*Chen et al., 2014*) and information theory (*Gong et al., 2013*).

94 While the diagnostics approach has helped establish a new philosophy and/or paradigm for hydrologic  
95 model evaluation, as a community, we continue to hold on to and rely too much on, deterministic,  
96 non-inclusive, measures of model performance. Point in case of these scoring functions is the so-called  
97 KG efficiency or KGE of *Gupta et al.* (2009) and improvements thereof (*Lamontagne et al., 2020*), which  
98 have gained rapid popularity among hydrologists as alternative to the infamous *Nash and Sutcliffe* (1970)  
99 efficiency, but are equally exclusive in how they quantify model performance using only single-valued  
100 simulated output of the quantities of interest. In fact, the confidence intervals of the KG efficiency are not  
101 defined in a formal mathematical sense preventing a rigorous statistical description of model parameter  
102 and simulation uncertainty (*Vrugt and de Oliveira, 2022*). A little more than a decade ago, *Guttorp*  
103 (2011) formulated a vision of how climate models should be evaluated against data (P. 820)

104 “... Climate models are difficult to compare to data. Often climatologists compute some  
105 summary statistic, such as global annual mean temperature, and compare climate models using  
106 observed (or rather estimated) forcings to the observed (or rather estimated) temperatures.  
107 However, it seems more appropriate to compare the distribution (over time and space) of  
108 climate model output to the corresponding distribution of observed data, as opposed to point  
109 estimates with or without confidence intervals.”

110 and this change from point to distributional evaluation finds support in hydrology through information-  
111 theoretic arguments (*Weijns et al. 2010a*, P. 2545)

112 “... models should preferably be explicitly probabilistic and calibrated to maximize the informa-  
113 tion they provide.”

114 and parallels sociocultural developments related to equity, inclusiveness and diversity. A recent report of  
115 McKinsey & Company titled *Diversity wins: How inclusion matters* (Hunt et al., 2020) concludes that  
116 companies with a more inclusive and greater workplace diversity outperform their competitors, increase  
117 employee engagement and achieve higher profits. These developments inspire a paradigm change in how  
118 we evaluate our models.

119 Distribution forecasts express diversity in the form of a probability distribution over future quantities  
120 or events (Dawid, 1984). Strictly speaking, the modeled outcomes in this paper are not forecasts as  
121 we use measured values of the exogenous variables. This, however, is inconsequential to the premise of  
122 this paper as methods discussed are equally applicable to simulation distributions. Such forecasts are  
123 routinely made in weather and climate prediction (Palmer, 2002; Gneiting et al., 2005), computational  
124 finance (Duffie and Pan, 1997), macroeconomics (Garratt et al., 2003; Granger, 2005) and hydrology  
125 (Welles et al., 2007; Thielen et al., 2008) but evade model evaluation by a need for ensemble simulation  
126 and lack of knowledge on how to properly evaluate a simulation distribution against measured data.  
127 Advances in computer hardware, software and theory have tremendously advanced our ability to quantify  
128 the uncertainty of hydrologic model output (Beven and Binley, 1992; Kuczera and Parent, 1998; Vrugt  
129 et al., 2003; Kavetski et al., 2006a; Schoups and Vrugt, 2010; Vrugt, 2016; Vrugt et al., 2022) but we  
130 typically discard distribution forecasts  $P$  of (quasi)-Bayesian methods in hydrologic model evaluation  
131 and quantify model performance using some set-valued mapping  $P \rightarrow T(P) \subseteq \Omega$  to the real line  $\mathbb{R}$  with  
132 the mean, likelihood-weighted mean or median prediction of the sample space  $\Omega$  as key examples. This  
133 point-valued mapping necessarily implies a huge loss of information about model performance. Exception  
134 to this are the works of Volpi et al. (2017) and Brunetti et al. (2017) who relied on the integrated or  
135 marginal likelihood for hydrologic model evaluation and selection.

136 Forecast verification is an active field of research in the climate, atmospheric and ocean sciences and is  
137 concerned with evaluating, verifying and determining the predictive power of prognostic model forecasts  
138 (Murphy and Katz, 1985; Storch and Zwiers, 1999; Jolliffe and Stephenson, 2011). Scoring rules have  
139 long been used to evaluate the accuracy of forecast probabilities after observing the occurrence, or  
140 nonoccurrence, of predicted events of dichotomous, categorical and continuous variables (Gneiting and  
141 Raftery, 2007).

142 **Definition 2.** *A scoring rule is any extended real-valued function  $S : \mathcal{P} \times \Omega \rightarrow \overline{\mathbb{R}} \equiv [-\infty, \infty]$  such that*

<sub>143</sub>  $S(P, \omega)$  is  $\mathcal{P}$ -quasi-integrable for all  $P \in \mathcal{P}$  and measures the reward (or loss) when the distribution  
<sub>144</sub> forecast  $P$  is issued and observation  $\omega \in \Omega$  materializes.

<sub>145</sub> Thus, a scoring rule  $S(P, \omega)$  measures the performance of a distribution forecast  $P$  in a single reward  
<sub>146</sub> (or loss) value and reduces to a scoring function  $s(y, \omega)$  for a point forecast. Most scoring rules are  
<sub>147</sub> real-valued, thus, take on values in  $\mathbb{R}$  with exceptions such as the ignorance score (*Roulston and Smith*,  
<sub>148</sub> 2002) or logarithmic rule (*Good*, 1952), which can attain scores of infinity and minus infinity, respectively,  
<sub>149</sub> and, thus, operate in  $\overline{\mathbb{R}}$ . The attractive statistical and information-theoretic properties of scoring rules  
<sub>150</sub> benefits ranking of likelihood functions (*Vrugt et al.*, 2022), hypothesis testing with watershed models  
<sub>151</sub> and, as we show in this paper, hydrologic model evaluation. All these are desirable qualities of scoring  
<sub>152</sub> rules given the plethora of hydrologic models used by researchers and practitioners (*Clark et al.*, 2008;  
<sub>153</sub> *Schoups et al.*, 2010; *Fenicia et al.*, 2011).

<sub>154</sub> *Weijis et al.* (2010a) presents a convincing example so as to why scoring rules such as the *Brier* (1950)  
<sub>155</sub> score and continuous ranked probability score (CRPS) of *Matheson and Winkler* (1976) should be used  
<sub>156</sub> for hydrologic model calibration and evaluation. Otherwise, model training is overly susceptible to  
<sub>157</sub> misinformation and/or incomplete (unfinished) learning. Despite their compelling plea, scoring rules  
<sub>158</sub> such as the CRPS have only found sporadic application in hydrology, usually for evaluating ensemble  
<sub>159</sub> forecast skill (*Vrugt et al.*, 2006; *Laio and Tamea*, 2007; *Girons Lopez et al.*, 2021). A simulation  
<sub>160</sub> distribution coalesces model responses across the (prior/posterior) parameter and/or input space and  
<sub>161</sub> contains information about model behavior, robustness, sensitivity and uncertainty that is not available  
<sub>162</sub> in single-valued model output. Thus, scoring function-based model evaluation strategies imply an  
<sub>163</sub> inherent loss of information about model functioning. This paper advocates the use of distribution-based  
<sub>164</sub> model evaluation and diagnostics. This paper is concerned with the basic question of how we should  
<sub>165</sub> evaluate distributions generated via simulation. We bring scoring rules to the attention of hydrologists  
<sub>166</sub> and demonstrate their power, usefulness and applicability to hydrologic model evaluation and model  
<sub>167</sub> diagnostics. We introduce *strictly proper* scoring rules for flow duration and recession curves and the  
<sub>168</sub> analysis of flood frequencies and extreme events. To understand the different scoring rules for dichotomous,  
<sub>169</sub> categorical and continuous variables, convey their relationship with information theory, explain the  
<sub>170</sub> importance of elicability and scoring rule propriety and sufficiency, we must review different concepts  
<sub>171</sub> from probability and information theory. Hopefully, our work inspires others to delve deeper into the

172 topic of scoring rules and seek the advantages of distribution-based model evaluation and diagnostics  
173 over the current practice of deterministic model evaluation.

174 The remainder of this paper is organized as follows. In Section 2, we formalize our mathematical/statistical  
175 treatment of probability and discuss our use of symbols and notation. Section 3 reviews the use of  
176 information theory, specifically relative entropy, applicable to ideal situations with known distribution  
177 functions of each verifying observation. Sections 4 - 7 discuss the more common and realistic situation in  
178 which we do not have knowledge of the distribution  $Q \in \mathcal{P}$  which materializes with the measurement  
179  $\omega \in \Omega$ . Section 4 illustrates the insufficiency of common methods used in the hydrologic literature for  
180 evaluating distribution forecasts. This is followed by Section 5, which relates information theory to  
181 scoring rules for distribution forecasts of categorical (discrete) variables and their extension in Section 6 to  
182 density forecasts of continuous variables. Section 7 is considered with diagnostic analysis and revisits the  
183 decomposition of *strictly proper* scoring rules into an uncertainty, sharpness (resolution) and reliability  
184 term. The different sections are permeated with case studies using conceptual watershed models and  
185 measurements of the rainfall-discharge transformation. The penultimate Section 8 presents a future  
186 outlook about the use of scoring rules in diagnostic model evaluation, Bayesian model selection and the  
187 prediction of extreme events. We introduce divergence scores for the flow duration curve and *Brutsaert*  
188 and *Nieber* (1977) recession analysis and present analytic expressions for CRPS and logarithmic scoring  
189 rules for the Pearson type III distribution used in flood frequency analysis. Section 9 concludes this  
190 paper with a summary of our main findings.

191 The topic of this paper warrants a strong statistical/information-theoretic treatment. But those discour-  
192 aged by our statistical treatment of this topic are directed to the hydrologic case studies of this paper  
193 and the **ScoringRules** toolbox in MATLAB. Proofs and computational details have been deferred to  
194 Appendices. Specifically, Appendix A reviews Gibbs' inequality. Appendices B, C, D, F, G, I and J  
195 derive closed-form expressions for the relative entropy and CRPS of certain parametric *forecast* and *true*  
196 distributions. Appendix E presents precipitation data. Finally, Appendix H details our computational  
197 and numerical implementation of the HYMOD, Hmodel and SAC-SMA conceptual watershed models.

<sup>198</sup> **2 Preliminaries**

<sup>199</sup> One of the major purposes of hydrologic modeling is to predict watershed behavior under future conditions.  
<sup>200</sup> We can shed much light on hydrologic theory, process knowledge, computational implementation, and  
<sup>201</sup> aleatory and epistemic uncertainty by formalizing what is involved in making such forecasts and by  
<sup>202</sup> assessing our methods on their empirical success at this task (e.g. *Dawid* 1984). Statistics provides  
<sup>203</sup> methods for the characterization of the uncertainty associated with future events or quantities. If  
<sup>204</sup> compelled by the interpretation of *Ramsey* (1926) and *de Finetti* (2017) that probability is a subjective  
<sup>205</sup> degree of belief, then the laws and theorems of probability theory will suffice to revise these subjective  
<sup>206</sup> probabilities (induction/learning) and express predictive uncertainty. Consequently, the probabilistic  
<sup>207</sup> forecasts in this paper are probability distributions over future events. We wish to quantify the statistical  
<sup>208</sup> consistency of the forecasts. This is a joint property of the forecasts and the events or values that  
<sup>209</sup> materialize.

<sup>210</sup> Before we proceed with a discussion of the underlying ideas, theory and principles, we first expose  
<sup>211</sup> our treatment of probability and explain our notation. We consider a *probabilistic forecast*,  $P$ , to be a  
<sup>212</sup> probability measure on the set of all possible outcomes of an experiment, the so-called sample space  
<sup>213</sup>  $\Omega$ . Let  $\Sigma$  be a nonempty collection of subsets of  $\Omega$  closed under complement, countable unions, and  
<sup>214</sup> countable intersections and let  $\mathcal{P}$  be a convex class of probability measures on  $(\Omega, \Sigma)$ . A *probabilistic*  
<sup>215</sup> *forecast* is a set function  $P \in \mathcal{P}$  from  $\Sigma$  to the real number line  $\mathbb{R} = (-\infty, \infty)$  which assigns probabilities  
<sup>216</sup>  $P \in [0, 1]$  to any subset  $\Sigma \subseteq \Omega$ , called an event, in a countably additive manner so that the entire sample  
<sup>217</sup> space has probability of one. Similarly, the *true forecast*  $Q \in \mathcal{P}$  assigns probabilities  $Q \in [0, 1]$  to all  
<sup>218</sup> events  $\Sigma \subseteq \Omega$ , with unit sum of all probabilities,  $Q(\Omega) = 1$ .

<sup>219</sup> The measure theoretic treatment of probability allows us to simultaneously treat discrete and continuous  
<sup>220</sup> probability distributions. In most of our examples, the sample space  $\Omega$  will consist of real numbers and  
<sup>221</sup>  $\Omega$  is synonymous to a random variable of interest, say, next day's peak discharge. The forecaster's task  
<sup>222</sup> is to quote a distribution  $P \in \mathcal{P}$  which characterizes the uncertainty of  $\Omega$ . Once the watershed has  
<sup>223</sup> revealed  $\omega \in \Omega$  the forecaster will obtain a reward  $S(P, \omega)$  depending on both the quoted distribution  
<sup>224</sup>  $P$  and the materialized value  $\omega$  of the peak discharge. The utility or reward function  $S$  is a so-called  
<sup>225</sup> *scoring rule* and plays a critical role in forecast evaluation. Note that a sign change reverses  $S$  to a  
<sup>226</sup> loss, cost or penalty function. A scoring rule is *strictly proper* if, regardless of the *true* distribution  $Q$  of

227 the peak discharge, the forecaster's expected reward is maximized if and only if he/she quotes  $P = Q$ .  
228 Such scoring rules encourage the forecaster to be honest and report the true set of probabilities. As in  
229 *McCarthy* (1956) we assume that the forecaster cannot influence the event predicted beyond the usual  
230 experimentation, data collection and modeling involved.

231 In recent decades, probabilistic forecasting methods have systematically found their way into the hydrologic  
232 literature for quantifying predictive uncertainty. The topic of forecast verification has not received such  
233 methodical treatment (but with exception of *Laio and Tamea* (2007)) and, as a result, methods for  
234 evaluating distribution forecasts, if used at all, are often used haphazardly. To expose the requirements  
235 of a meaningful evaluation of distribution forecast and establish the relationship with information theory  
236 we have to be comprehensive in our statistical treatment. If deemed appropriate, we spare the main  
237 text from lengthy mathematical derivations and defer such technicalities and the description of data,  
238 models and computational procedures to appendices. A lowercase italic font ( $a$ ) is used for scalars, a  
239 lowercase bold font ( $\mathbf{a}$ ) for vectors and an uppercase bold font ( $\mathbf{A}$ ) for matrices. Moreover, the standard  
240 mathematical convention of functions, say  $f(x)$  and  $g(x)$ , will be used throughout. The symbol  $\omega$  is  
241 used for verifying measurement; thus, we write  $\omega_1, \dots, \omega_n$  for the time series of materialized outcomes.  
242 Statistical distributions are designated common symbols. If  $x$  is an observable random quantity with a  
243 normal distribution of mean  $\mu \in \mathbb{R}$  and variance  $\sigma^2 > 0$ , we write  $x \sim \mathcal{N}(\mu, \sigma^2)$ ; if the distribution of  $x$  is  
244 binomial with number of trials  $n \in \mathbb{N}_+$  and probability of success  $p \in [0, 1]$  we write,  $x \sim \mathcal{B}(n, p)$  and use  
245  $x \sim \mathcal{U}(a, b)$  for the continuous uniform distribution on the closed-interval  $[a, b]$ , where  $a, b \in \mathbb{R}$  and  $a < b$ .  
246 We designate a probability density function (PDF) with a lowercase  $f$  and a cumulative distribution  
247 function (CDF) with an uppercase  $F$ . Thus,  $f_{\mathcal{N}}(x, \mu, \sigma^2)$  and  $F_{\mathcal{N}}(x, \mu, \sigma^2)$  signify the PDF and CDF  
248 of the normal distribution, respectively. We use the vertical bar " $|$ " for conditional expectation. Thus,  
249  $p(x|\boldsymbol{\omega})$  is the conditional PDF of  $x$  given the  $n$ -vector of verifying measurements  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_n)^\top$   
250 with  $p(x|\boldsymbol{\omega}) \geq 0$  and  $\int_{\Omega} p(x|\boldsymbol{\omega}) dx = 1$ .

### 251 3 Relative entropy

252 Let us assume that we have exact knowledge of the distribution  $Q$  of the measurement,  $\omega$ , that will  
253 materialize at some future time. This situation, albeit uncommon, is a logical starting point for our

discussion. In mathematical statistics, the *Kullback and Leibler* (1951) divergence,  $d_{\text{KL}}(Q, P)$ , also known as relative entropy or  $I$ -divergence (*Csiszar*, 1975), measures in a single value the distance between the *forecast distribution*  $P$  and a reference or *true distribution*  $Q$  (*Kullback and Leibler*, 1951; *Kullback*, 1959). Divergence is a physical measure of information gain in communication theory and machine learning. For the time being we adapt the common practice in information theory and precede the *probabilistic forecast*  $P$  with the *true distribution*  $Q$ . These two arguments are swapped later in the context of scoring rules.

### 3.1 Continuous random variables

For distributions  $Q$  and  $P$  of a continuous random variable in sample space  $\Omega$ , the relative entropy from  $P$  to  $Q$  is defined as follows (*Jaynes*, 1963)

$$d_{\text{KL}}(Q, P) = \mathbb{E}_Q \left[ \log_b \left( \frac{Q(x)}{P(x)} \right) \right] = \int_{x \in \Omega} Q(x) \log_b \left( \frac{Q(x)}{P(x)} \right) dx, \quad (1)$$

where  $Q(x)$  and  $P(x)$  are the probabilities of  $Q$  and  $P$  evaluated at the event  $x \in \Omega$  and  $Q(\Omega) = 1$  and  $P(\Omega) = 1$ . In applications,  $Q$  typically signifies the *true distribution* of data, observations, or possibly, some exactly defined theoretical distribution, whereas  $P$  is an approximation thereof obtained from paper-and-pencil calculation, computer modeling and/or other quantitative means. Note that our assignment of the symbols  $Q$  and  $P$  to the true and forecast distribution, respectively, is reversed to common practice in information theory but consistent with the statistical literature on forecast evaluation.

The symbol  $b$  is used for the base of the logarithm. Common values of  $b$  are 2,  $e = 2.7182818\dots$  (Euler's number) and 10 and give units of the (relative) entropy in bits (or shannons), nats and hartleys (also referred to as dits or bans), respectively. In what follows we do not affix the base  $b$  of the logarithm and assume units of bits in our colloquial references to entropy.

The relative entropy  $d_{\text{KL}}(Q, P)$  is defined only if the ratio  $Q(x)/P(x)$  of the two probability measures, the so-called Radon-Nikodym derivative,  $dQ/dP$ , exists. This means that there does not exist an event  $x \in \Omega$  for which  $Q(x) > 0$  and  $P(x) = 0$ , otherwise we must divide by zero. As  $\log_b(a/b) = \log_b(a) - \log_b(b)$ , the familiar information-theoretic expression for the relative entropy is

$$\begin{aligned} d_{\text{KL}}(Q, P) &= \int_{x \in \Omega} \left( Q(x) \log_b(Q(x)) - Q(x) \log_b(P(x)) \right) dx \\ &= \underbrace{\int_{x \in \Omega} Q(x) \log_b(Q(x)) dx}_{-\mathbb{H}(Q)} - \underbrace{\int_{x \in \Omega} Q(x) \log_b(P(x)) dx}_{-\mathbb{H}(Q, P)} \end{aligned}$$

$$= \mathbb{H}(Q, P) - \mathbb{H}(Q). \quad (2)$$

where  $\mathbb{H}(Q, P)$  is the so-called cross-entropy between the *true distribution*  $Q$  and the *probabilistic forecast*  $P$  and  $\mathbb{H}(Q)$  is the Shannon entropy of the *true distribution*  $Q$  itself (Shannon, 1948a,b). The cross-entropy measures the number of bits required to represent or transmit an average event from distribution  $Q$  compared to distribution  $P$ . If  $Q \neq P$  the cross-entropy  $\mathbb{H}(Q, P)$  will always exceed the entropy  $\mathbb{H}(Q)$  and  $d_{\text{KL}}(Q, P) > 0$ . This is known as Gibbs' inequality, a common proof of which is reiterated in Appendix A. If  $P = Q$  and our *distribution forecast*  $P$  matches exactly the *true distribution*  $Q$  then  $\mathbb{H}(Q, P)$  will equal  $\mathbb{H}(Q)$  and  $d_{\text{KL}}(Q, P)$  is zero. Thus,  $d_{\text{KL}}(Q, P) = 0$  if and only if  $P = Q$ . In conclusion, the closer the value of the relative entropy  $d_{\text{KL}}(Q, P)$  to zero, the more similar  $Q$  and  $P$  will be.

If  $Q$  and  $P$  are strictly continuous on  $\mathcal{P}$  and follow a known statistical distribution then it may be possible to derive an analytic expression for their distance,  $d_{\text{KL}}(Q, P)$  (see e.g. Bouhlel and Dziri 2019). Suppose, for example, that the *true probability distribution*  $Q$  of some quantity of interest  $\Omega$  equals the uniform distribution,  $Q = \mathcal{U}(a, b)$  on the closed interval  $\Omega = [a, b]$  with PDF,  $f_{\mathcal{U}}(\omega, a, b) = 1/(b - a)$  (see Figure 1). The distribution *forecast*  $P$  of  $\Omega$  equals a symmetric triangular distribution,  $P = \mathcal{T}(a, b)$ , with midpoint  $c = (a + b)/2$  and PDF

$$f_{\mathcal{T}}(x, a, b) = \frac{2(b - a) - 2|a + b - 2x|}{(b - a)^2}, \quad (3)$$

where  $|\cdot|$  denotes the absolute value. We can enter the analytic expressions of the PDFs of  $Q$  and  $P$  into the integral of Equation (1) to yield  $d_{\text{KL}}(Q, P) = 1 - \log_b(2)$  in Appendix B. Note that  $d_{\text{KL}}(Q, P)$  does not depend on the width of the interval of  $\Omega$ . Interestingly, for  $b = 2$  the relative entropy  $d_{\text{KL}}(Q, P) = 0$  even though  $Q \neq P$ . This artifact is easily resolved with a temporary change to the base of the logarithm to yield,  $d_{\text{KL}}(Q, P) = (1 - \log_e(b))/\log_e(b)$ . If we then admit  $b = 2$  we yield  $d_{\text{KL}}(Q, P) = 0.4427$  bits. If we swap the arguments  $Q$  and  $P$  in our derivation, we yield the relative entropy from  $Q$  to  $P$  instead. In Appendix B, we derive that  $d_{\text{KL}}(P, Q) = \log_b(2) - \frac{1}{2}$  or  $d_{\text{KL}}(P, Q) = 0.2787$  bits.

The analytic exercise demonstrates that the relative entropy does not satisfy the symmetry axiom of a metric  $d : \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}_+$  in a metric space  $\mathcal{M}$ . Indeed, the relative entropy  $d_{\text{KL}}(Q, P)$  from  $P$  to  $Q$  does not equal its counterpart  $d_{\text{KL}}(P, Q)$  from  $Q$  to  $P$ . But this is not all. The relative entropy also does not satisfy the triangle inequality,  $d_{\text{KL}}(Q, P) \leq d_{\text{KL}}(Q, R) + d_{\text{KL}}(R, P)$ , which we will confirm later with an analytic example. This fourth and last axiom of a metric  $d$  in  $\mathcal{M}$  states that the distance from  $Q$  to

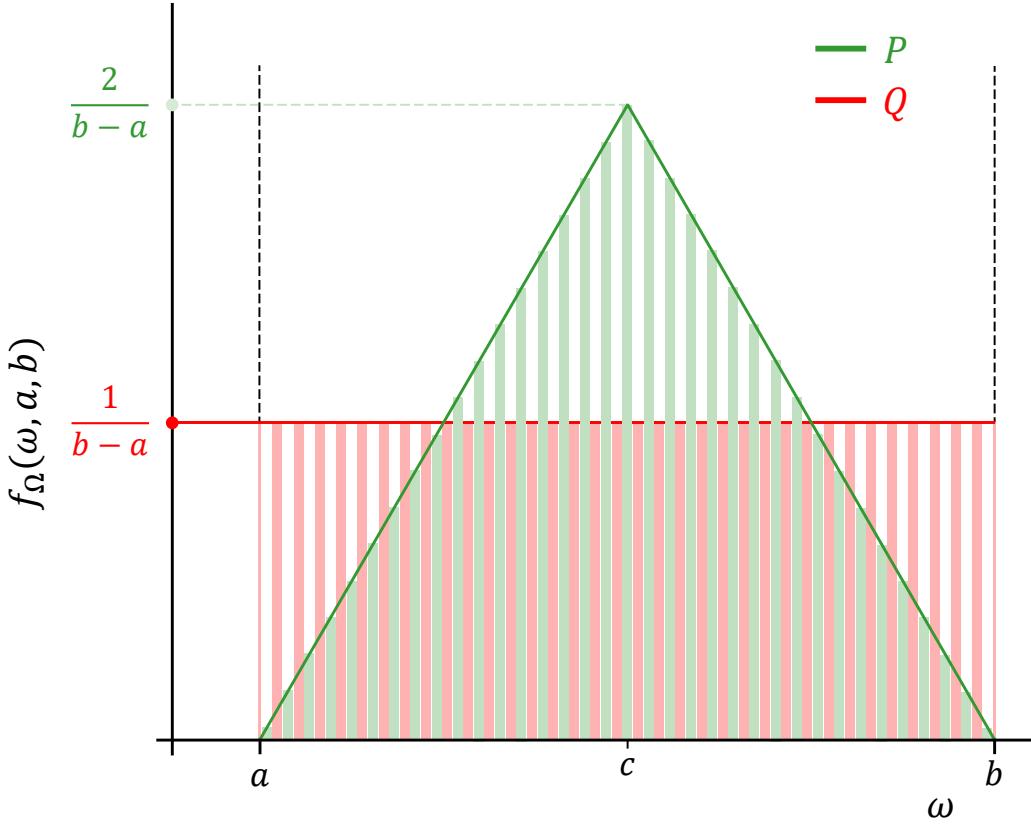


Figure 1: Probability density function of the uniform true distribution  $Q = \mathcal{U}(a, b)$  and symmetric triangular forecast distribution  $P = \mathcal{T}(a, b)$  of quantity  $\Omega = [a, b]$ .

310  $P$  with a detour through  $R$  cannot be less than the geodesic distance (shortest path) between  $Q$  and  
 311  $P$ . This implies that relative entropy is not a metric or distance function in an Euclidean space with  
 312 its usual physical or metaphorical notion of length but rather should be thought of as a divergence of  
 313 two distributions. This statistical distance is commonly referred to in the statistical literature as the  
 314 *Kullback and Leibler* (1951) divergence, hence our use of the subscript KL in  $d_{\text{KL}}(Q, P)$ . Furthermore,  
 315 the fact that  $d_{\text{KL}}(Q, P)$  is strictly nonnegative and zero only when  $P = Q$  suggests its interpretation  
 316 as a so-called *Bregman* (1967) *distance* or *divergence* of  $Q$  and  $P$ . We point forward to Figure 5 for a  
 317 graphical exposition of the Bregman divergence but refrain from further explanation until we get to the  
 318 topic of scoring rules in Section 5.

319 Bregman divergences are of paramount importance in scientific forecast evaluation as they ensure *strict*  
 320 *propriety* of scoring rules and their associated divergence functions. Thus,  $d_{\text{KL}}(Q, P) > 0$  if  $P \neq Q$  and  
 321  $d_{\text{KL}}(Q, P) = 0$  if and only if  $P = Q$ . Propriety, “*... the state or quality of conforming to conventionally*  
 322 *accepted standards of behavior or morals and/or ... conformity with what is required by a rule, principle,*  
 323 *etc.*” follows directly from Jensen’s inequality (the convex transformation of a mean is less than or equal

324 to the mean applied after convex transformation) and incentivizes a forecaster to issue  $P = Q$  rather than  
 325 any  $P \neq Q$ . In other words, only *strictly proper* scoring rules will lead us to the *true forecast distribution*,  
 326  $Q \in \mathcal{P}$ . The minimum divergence,  $d(P, Q) = 0$  is achieved when  $P = Q$ , and this minimum is unique.  
 327 This encourages a forecaster to be honest and to volunteer his or her true beliefs. *Proper* scoring rules  
 328 also yield a minimum at  $P = Q$  but this minimum may not be unique. A symmetrized variant of the  
 329 *Kullback and Leibler* (1951) divergence

$$d_J(P, Q) = d_{KL}(Q, P) + d_{KL}(P, Q) = \int_{x \in \Omega} (Q(x) - P(x)) \left( \frac{Q(x)}{P(x)} \right) dx \quad (4)$$

330 is also known as *J*-divergence in honor of Sir Harold Jeffreys (*Jeffreys*, 1946) and used widely in pattern  
 331 recognition and computer vision applications (*Chang et al.*, 2009; *Zheng and You*, 2013). Divergences are  
 332 sometimes called divergence functions or discrepancy functions, or validation metrics (*Liu et al.*, 2011),  
 333 even though our example has shown that they may not necessarily satisfy the requirements of a metric  
 334 in mathematical sense (*Thorarinsdottir et al.*, 2013).

335 If the *true distribution*  $Q$  and *probabilistic forecast*  $P$  are multivariate normal in  $\mathbb{R}^b$  with means,  
 336  $\mu_Q, \mu_P \in \mathbb{R}^{b \times 1}$ , and non-singular  $b \times b$  covariance matrices  $\Sigma_Q$  and  $\Sigma_P$ , respectively, then the relative  
 337 entropy in units of nats becomes (see Appendix C)

$$d_{KL}(\mathcal{N}_b(\mu_Q, \Sigma_Q), \mathcal{N}_b(\mu_P, \Sigma_P)) = \frac{1}{2} [\log_e(|\Sigma_Q^{-1} \Sigma_P|) - b + \text{tr}(\Sigma_P^{-1} \Sigma_Q) + (\mu_Q - \mu_P)^\top \Sigma_P^{-1} (\mu_Q - \mu_P)] \quad (5)$$

338 where  $|\cdot|$  is the determinant operator, the symbol  $\top$  denotes transpose and the trace function,  $\text{tr}(\mathbf{A})$ ,  
 339 returns the sum of the elements on the main diagonal of the  $b \times b$  matrix,  $\mathbf{A} = \Sigma_1^{-1} \Sigma_0$ . Equation (5)  
 340 is also known as the *Dawid and Sebastiani* (1999) divergence,  $d_{DDS}(P, Q)$ , more of which later in the  
 341 context of multivariate scoring rules. Note that the arguments in  $d_{DDS}(P, Q)$  have reversed to satisfy  
 342 convention used in the statistical literature. This analytic expression of the KL-divergence is *strictly*  
 343 *proper* only for normal probability measures, which are uniquely characterized by their respective means  
 344 and covariance matrices. If  $Q$  and  $P$  are univariate normal then Equation (5) reduces to

$$d_{KL}(\mathcal{N}_1(\mu_Q, \sigma_Q^2), \mathcal{N}_1(\mu_P, \sigma_P^2)) = \frac{1}{2} \log_e \left( \frac{\sigma_P^2}{\sigma_Q^2} \right) + \frac{\sigma_Q^2 + (\mu_Q - \mu_P)^2 - \sigma_P^2}{2\sigma_P^2}. \quad (6)$$

345 The relative entropy  $d_{KL}(Q, P)$  is not defined when  $P(x) = 0$  and  $Q(x) > 0$ . To illustrate this, suppose  
 346  $Q = \mathcal{N}(\mu_Q, \sigma_Q^2)$  and  $P = \mathcal{U}(a_P, b_P)$  on a finite sample space  $\Omega = [a, b]$ . In Appendix D we derive an

351 analytic expression for the KL-divergence of  $Q$  from  $P$

$$\begin{aligned} d_{\text{KL}}(\mathcal{N}(\mu_Q, \sigma_Q^2), \mathcal{U}(a_P, b_P)) &= \mathbb{H}(\mathcal{N}(\mu_Q, \sigma_Q^2), \mathcal{U}(a_P, b_P)) - \mathbb{H}(\mathcal{N}(\mu_Q, \sigma_Q^2)) \\ &= \log_e(b_P - a_P) - \frac{1}{2} \log_e(2e\pi\sigma_Q^2), \end{aligned} \quad (7)$$

352

353

355 and may go to infinity with support of  $P$  on the extended real line,  $\overline{\mathbb{R}}$ . Note that if we link  $Q$  and  $P$  using  
356  $\sigma = (b_P - a_P)/\nu$  with  $\nu \in \mathbb{R}_+$  then the relative entropy simplifies to  $d_{\text{KL}}(Q, P) = \log_e(\nu) - \frac{1}{2} \log_e(2e\pi)$ .

357 Now suppose that we have three distributions,  $Q = \mathcal{N}(\mu_Q, \sigma_Q^2)$ ,  $P = \mathcal{U}(a_P, b_P)$  and  $R = \mathcal{N}(\mu_R, \sigma_R^2)$ .  
358 Equations (6) and (7) will help demonstrate that the relative entropy does not honor the triangle  
359 inequality,  $d_{\text{KL}}(Q, P) \leq d_{\text{KL}}(Q, R) + d_{\text{KL}}(R, P)$ . Indeed, we yield

$$\begin{aligned} d_{\text{KL}}(\mathcal{N}(\mu_Q, \sigma_Q^2), \mathcal{U}(a_P, b_P)) &\leq d_{\text{KL}}(\mathcal{N}(\mu_Q, \sigma_Q^2), \mathcal{N}(\mu_R, \sigma_R^2)) + d_{\text{KL}}(\mathcal{N}(\mu_R, \sigma_R^2), \mathcal{U}(a_P, b_P)) \\ \log_e(b_P - a_P) - \frac{1}{2} \log_e(2e\pi\sigma_Q^2) &\leq \frac{1}{2} \log_e\left(\frac{\sigma_R^2}{\sigma_Q^2}\right) + \frac{\sigma_Q^2 + (\mu_Q - \mu_R)^2 - \sigma_R^2}{2\sigma_R^2} + \log_e(b_P - a_P) - \frac{1}{2} \log_e(2e\pi\sigma_R^2) \\ -\frac{1}{2} \log_e(\sigma_Q^2) &\leq \frac{1}{2} \log_e\left(\frac{\sigma_R^2}{\sigma_Q^2}\right) + \frac{\sigma_Q^2 + (\mu_Q - \mu_R)^2 - \sigma_R^2}{2\sigma_R^2} - \frac{1}{2} \log_e(\sigma_R^2) \\ \implies \sigma_Q^2 + (\mu_Q - \mu_R)^2 - \sigma_R^2 &\geq 0. \end{aligned} \quad (8)$$

360

361

362

363

365 The trivial example,  $\sigma_Q^2 < \sigma_R^2$  and  $\mu_Q = \mu_R$ , violates the triangle inequality. To convey the fundamental  
366 asymmetry in the relation between  $Q$  and  $P$  it is common to refer to  $d_{\text{KL}}(Q, P)$  as the relative entropy  
367 of  $Q$  with respect to  $P$  or the information gain from  $Q$  over  $P$ .

## 368 3.2 Discrete random variables

369 For discrete probability distributions  $Q$  and  $P$  the sample space,  $\Omega = \{\omega_1, \dots, \omega_m\}$  consists of a finite  
370 number  $m$  of mutually exclusive and collectively exhaustive events,  $\omega$ , and a probabilistic forecast is a  
371 probability vector  $\mathbf{p} = (p_1, \dots, p_m)^\top$  defined on the convex class  $\mathcal{P} = \mathcal{P}_m$

$$\mathcal{P}_m = \{\mathbf{p} = (p_1, \dots, p_m)^\top : p_1 + \dots + p_m = 1 \text{ and } p_k \geq 0 \text{ for all } k\}. \quad (9)$$

373 where  $\top$  denotes transpose. It is further assumed that the vector  $\mathbf{q} = (q_1, \dots, q_m)^\top$  reports the *true*  
374 *probabilities* of the  $m$  events,  $\{\mathbf{q} \in \mathbb{R}_+^{m \times 1} : \mathbf{1}^\top \mathbf{q} = 1\}$ , where  $\mathbf{1}_m$  is a  $m \times 1$  vector of ones.

375 For discrete probability distributions  $P$  and  $Q$  on sample space  $\Omega$ , Equation (2) reduces to

$$d_{\text{KL}}(\mathbf{q}, \mathbf{p}) = \mathbb{H}(\mathbf{q}, \mathbf{p}) - \mathbb{H}(\mathbf{q}), \quad (10)$$

376

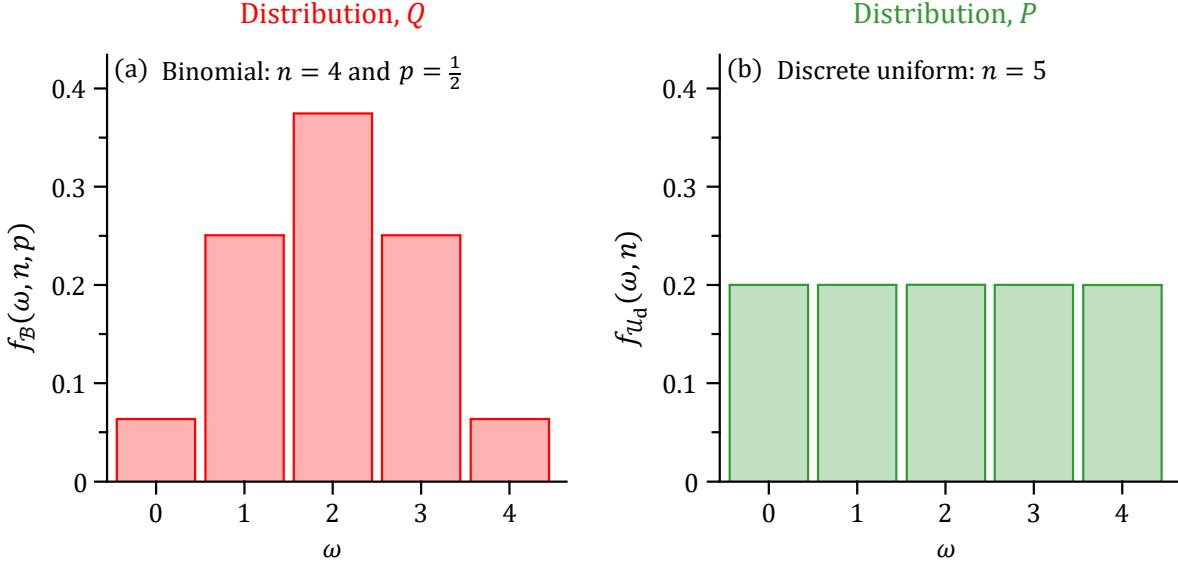


Figure 2: Illustration of the computation of the relative entropy  $d_{\text{KL}}(Q, P)$  on a sample space  $\Omega = \{0, 1, 2, 3, 4\}$ : (a) *true distribution*,  $Q = \mathcal{B}(n, p)$ , and (b) *forecast distribution*,  $P = \mathcal{U}_d(n)$ . According to data,  $\Omega$  follows a binomial distribution with  $n = 4$ ,  $p = \frac{1}{2}$  and PMF,  $f_B(\omega, n, p) = c(n, \omega)p^\omega(1-p)^{n-\omega}$ , where  $c(a, b) = a!/b!(a-b)!$  denotes the binomial coefficient and ! is the factorial function. Theory predicts a discrete uniform distribution for  $\omega$  with equal density,  $f_{\mathcal{U}_d}(\omega, n) = 1/n$ , for all  $n = 5$  values.

377 and the integral of the relative entropy from  $P$  to  $Q$  becomes a nonnegative sum

$$378 \quad d_{\text{KL}}(\mathbf{q}, \mathbf{p}) = \sum_{k=1}^m q_k \log_b \left( \frac{q_k}{p_k} \right), \quad (11)$$

379 which is equivalent to

$$380 \quad d_{\text{KL}}(\mathbf{q}, \mathbf{p}) = - \sum_{k=1}^m q_k \log_b \left( \frac{p_k}{q_k} \right). \quad (12)$$

381 To illustrate the computation and interpretation of the KL-divergence, please consider the probability  
382 mass functions (PMFs) of  $Q$  and  $P$  depicted in Figure 2. The *true distribution*  $Q$  of quantity  $\Omega$  is a  
383 binomial distribution with  $n = 4$  and  $p = 0.5$  and the *distribution forecast*  $P$  is discrete uniform with  
384 equal density for  $\omega = (0, 1, \dots, 4)$ . The relative entropy,  $d_{\text{KL}}(\mathbf{q}, \mathbf{p})$ , may now be computed

$$\begin{aligned} 385 \quad d_{\text{KL}}(\mathbf{q}, \mathbf{p}) &= \sum_{\omega=0}^4 f_B(\omega, 4, \frac{1}{2}) \log_b \left( \frac{f_B(\omega, 4, \frac{1}{2})}{f_{\mathcal{U}_d}(\omega, 5)} \right) \\ 386 \quad &= \frac{1}{16} \log_b \left( \frac{\frac{1}{16}}{\frac{1}{5}} \right) + \frac{1}{4} \log_b \left( \frac{\frac{1}{4}}{\frac{1}{5}} \right) + \frac{6}{16} \log_b \left( \frac{\frac{6}{16}}{\frac{1}{5}} \right) + \frac{1}{4} \log_b \left( \frac{\frac{1}{4}}{\frac{1}{5}} \right) + \frac{1}{16} \log_b \left( \frac{\frac{1}{16}}{\frac{1}{5}} \right) \\ 387 \quad &= \frac{1}{8} \log_b \left( \frac{5}{16} \right) + \frac{1}{2} \log_b \left( \frac{5}{4} \right) + \frac{3}{8} \log_b \left( \frac{30}{16} \right), \end{aligned} \quad (13)$$

389 which with base of the logarithmic function equal to two leads to  $d_{\text{KL}}(\mathbf{q}, \mathbf{p}) \approx 0.2913$  bits. If we divide  
390 the so-obtained value of  $d_{\text{KL}}(\mathbf{q}, \mathbf{p})$  by  $\log_2(e)$  then we yield the KL-divergence in units of nats. Note

391 that if  $q_k = 0$  for some  $\omega_k \in \Omega$ , the summand,  $q_k \log_b(q_k)$ , is set to zero in accordance with the limit,  
 392  $\lim_{q \downarrow 0} q \log_b(q) = 0$ . If we swap the discrete distributions of  $Q$  and  $P$  we yield the so-called reverse  
 393 KL-divergence

$$\begin{aligned} 394 \quad d_{\text{KL}}(\mathbf{p}, \mathbf{q}) &= \sum_{\omega=0}^4 f_{U_d}(\omega, 5) \log_b \left( \frac{f_{U_d}(\omega, 5)}{f_B(\omega, 4, \frac{1}{2})} \right) \\ 395 \quad &= \frac{1}{5} \log_b \left( \frac{\frac{1}{5}}{\frac{1}{16}} \right) + \frac{1}{5} \log_b \left( \frac{\frac{1}{5}}{\frac{1}{4}} \right) + \frac{1}{5} \log_b \left( \frac{\frac{1}{5}}{\frac{6}{16}} \right) + \frac{1}{5} \log_b \left( \frac{\frac{1}{5}}{\frac{1}{4}} \right) + \frac{1}{5} \log_b \left( \frac{\frac{1}{5}}{\frac{1}{16}} \right) \\ 396 \quad &= \frac{2}{5} \log_b \left( \frac{16}{5} \right) + \frac{2}{5} \log_b \left( \frac{4}{5} \right) + \frac{1}{5} \log_b \left( \frac{16}{30} \right), \\ 397 \end{aligned} \quad (14)$$

398 which is equal to  $d_{\text{KL}}(\mathbf{p}, \mathbf{q}) \approx 0.3611$  bits. This again confirms that  $d_{\text{KL}}(\mathbf{q}, \mathbf{p})$  is not a metric but rather  
 399 a (Bregman) divergence, more of which later.

400 Expressed in the language of Bayesian inference,  $d_{\text{KL}}(Q, P)$  is a measure of the information gained by  
 401 revising one's beliefs from the prior probability distribution  $P$  to the posterior probability distribution  $Q$ .  
 402 In other words, it is the amount of information lost when  $P$  is used to approximate  $Q$  (*Burnham and*  
 403 *Anderson, 2002*). We used the KL-divergence in our earlier work to quantify the gain of information when  
 404 moving from the prior to the posterior parameter distribution (*Scharnagl et al., 2010*). For example, if  
 405  $\Theta = (\theta_1, \dots, \theta_d)^\top$  equals a  $d$ -vector of unknown coefficients of some mathematical model or hypothesis,  $\mathcal{H}$ ,  
 406 with prior distribution,  $p(\Theta|\mathcal{H})$  on the parameter (sample) space  $\Theta \subseteq \mathbb{R}^d$ . When new data,  $\mathcal{D}$ , become  
 407 available, we can update  $p(\Theta|\mathcal{H})$  to a posterior distribution,  $p(\Theta|\mathcal{D}, \mathcal{H})$ , using *Bayes* (1763) theorem

$$408 \quad p(\Theta|\mathcal{D}, \mathcal{H}) = \frac{p(\mathcal{D}|\Theta, \mathcal{H})p(\Theta|\mathcal{H})}{p(\mathcal{D}|\mathcal{H})}. \quad (15)$$

409 where  $p(\mathcal{D}|\mathcal{H})$  is a normalization constant and  $L(\Theta|\mathcal{D}, \mathcal{H}) \equiv p(\mathcal{D}|\Theta, \mathcal{H})$  denotes the likelihood. The  
 410 entropy of the posterior distribution

$$411 \quad \mathbb{H}(p(\Theta|\mathcal{D}, \mathcal{H})) = - \sum_{\Theta \in \Theta} p(\Theta|\mathcal{D}, \mathcal{H}) \log_b(p(\Theta|\mathcal{D}, \mathcal{H})), \quad (16)$$

412 may be less than, equal to, or greater than the entropy  $\mathbb{H}(p(\Theta|\mathcal{H}))$  of the prior distribution. The relative  
 413 entropy

$$414 \quad d_{\text{KL}}(p(\Theta|\mathcal{D}, \mathcal{H}), p(\Theta|\mathcal{H})) = \mathbb{H}(p(\Theta|\mathcal{D}, \mathcal{H}), p(\Theta|\mathcal{H})) - \mathbb{H}(p(\Theta|\mathcal{D}, \mathcal{H})), \quad (17)$$

415 measures the added message length (bits) that an original code with the prior distribution  $p(\Theta|\mathcal{H})$   
 416 would require relative to a new code based on the posterior distribution  $p(\Theta|\mathcal{D}, \mathcal{H})$ . Or in the words of  
 417 *Cover and Thomas* (2006) (Chapter 2, Page 19) “... if we knew the true distribution  $Q$  of the random

418 variable, we could construct a code with average description length  $\mathbb{H}(Q)$ . If, instead, we used the  
 419 code for a distribution  $P$ , we would need  $\mathbb{H}(Q) + d_{\text{KL}}(Q, P)$  bits on the average to describe the random  
 420 variable”, where  $Q$  and  $P$  correspond to the posterior and prior distributions, respectively. Thus,  
 421  $d_{\text{KL}}(p(\boldsymbol{\theta}|\mathcal{D}, \mathcal{H}), p(\boldsymbol{\theta}|\mathcal{H}))$ , is synonymous with the information gain about  $\boldsymbol{\Theta}$  that has been learned by the  
 422 new data  $\mathcal{D}$ . Note, that  $p(\boldsymbol{\theta}|\mathcal{D}, \mathcal{H})$  (and  $p(\boldsymbol{\theta}|\mathcal{H})$  for that matter) must equal a probability measure on  $\boldsymbol{\Theta}$   
 423 and, thus,  $\sum_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} p(\boldsymbol{\theta}|\mathcal{D}, \mathcal{H}) = 1$ . Posterior distributions derived from MCMC methods lend itself well to  
 424 computation of the Shannon entropy via Equation (16). If regularity conditions are satisfied then the  
 425 Markov chain will visit each event  $\boldsymbol{\theta}_k \in \boldsymbol{\Theta}$  with frequency  $n_k$  proportional to its so-called Boltzmann  
 426 weight,  $Q(\boldsymbol{\theta}_k) \sim \exp(-U(\boldsymbol{\theta}_k)/k_B T)$ , where  $U(\boldsymbol{\theta}_k)$  is the potential energy at event  $\boldsymbol{\theta}_k$ ,  $T$  is the absolute  
 427 temperature and  $k_B = 1.380649 \times 10^{-23}$  JK $^{-1}$  signifies the Boltzmann constant. The relative frequency  
 428 of each of  $m$  posterior samples (events or micro-states) visited by the Markov chain is given by the  
 429 Boltzmann distribution

$$Q(\boldsymbol{\theta}_k) = \frac{n_k}{m} = \frac{\exp(-U(\boldsymbol{\theta}_k)/k_B T)}{\sum_{i=1}^m \exp(-U(\boldsymbol{\theta}_i)/k_B T)}, \quad (18)$$

430 and defines a probability measure  $Q \in \mathcal{P}$  with  $Q(\boldsymbol{\Theta}) = 1$ . For arbitrary distributions  $Q$  and  $P$  in  $\mathbb{R}^d$  we  
 431 must evaluate Equation (11) at a collection of points (events),  $\mathbf{x} \in \Omega$ . This is equivalent to numerical  
 432 integration and different approaches exist to do so efficiently including dimension-adaptive quadrature  
 433 rules (Smolyak, 1963; Gerstner and Griebel, 1998, 2003; Jakeman and Narayan, 2018), quasi-random  
 434 sequences (Halton, 1960; Niederreiter, 1992; Sobol' and Shukhman, 1995) or pseudo-random sampling  
 435 (Hammersley and Handscomb, 1960; Volpi et al., 2017).

## 4 Insufficient scoring rules for density forecasts

437 In most practical situations we do not have knowledge of the underlying distribution  $Q \in \mathcal{P}$  which  
 438 materializes with the measurement  $\omega \in \Omega$ . Then, we must find ways to evaluate the distribution forecast,  
 439  $P \in \mathcal{P}$ , using only a single verifying observation,  $\omega_{st}$ , at given space and time coordinates. The data  
 440 may arrive *en bloc* in simulation mode or in natural (sequential) order in a forecasting problem. Then, it  
 441 is the forecaster's task, at any time, to produce a distribution forecast for the next observation. The  
 442 success at this task can be judged by using methods from probability forecasting. Before we address the  
 443 intriguing topic of scoring rules, we must first review past developments that led to current perspectives.

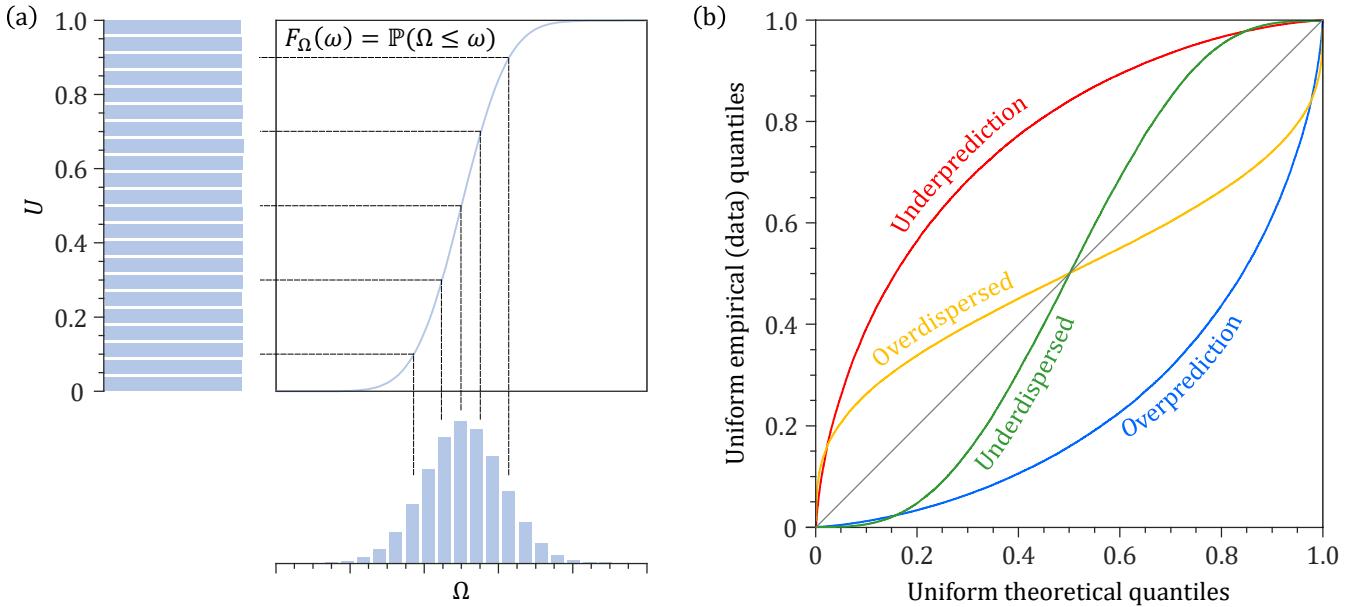


Figure 3: (a) Illustration of the probability integral transform and (b) interpretation of the so-called quantile-quantile plot (adapted from *Laio and Tamea (2007)*).

<sup>445</sup> This is a necessary step into understanding the strengths and weaknesses of current metrics used in  
<sup>446</sup> hydrology and the need for scoring rules.

<sup>447</sup> The probability integral transform of *Dawid* (1984) is one of the earliest methods for evaluating the  
<sup>448</sup> statistical coherency and association between a time series of forecast distributions  $P_1, \dots, P_n$  and observed  
<sup>449</sup> outcomes  $\omega_1, \dots, \omega_n$ . Integral transforms have a long history dating back to at least the *Rosenblatt*  
<sup>450</sup> (1952) transformation and turn a vector of dependent random variables into a vector of independent  
<sup>451</sup> uniform distributed values (see Figure 3a). Let  $\Omega$  be a real-valued continuous random variable with  
<sup>452</sup> CDF,  $F_\Omega(\omega) = \mathbb{P}(\Omega \leq \omega)$ . Then, random variable  $U = F_\Omega(\omega)$  has a standard uniform distribution. Thus,  
<sup>453</sup> if  $\omega_1, \dots, \omega_n$  are samples of  $\Omega$  (dependent or not) then the values  $u_i = F_\Omega(\omega_i)$ ;  $i = (1, \dots, n)$  of  $U$  will  
<sup>454</sup> be uniformly distributed on the unit interval. Hence, the probability integral transform reduces the  
<sup>455</sup> assessment of  $F_\Omega$  to the question whether the sequence of  $u$ 's behaves as a random sample of  $\mathcal{U}[0, 1]$ . Fig.  
<sup>456</sup> 3b illustrates the consequences of using an incorrect distribution for  $\Omega$  on the relationship between the  
<sup>457</sup> theoretical quantiles of the standard uniform distribution and the quantiles of the empirical distribution  
<sup>458</sup> function of the sampled data,  $\omega_1, \dots, \omega_n$ . This so-called predictive quantile-quantile (Q-Q) plot (*Dawid*,  
<sup>459</sup> 1984; *Casella and Berger*, 2002) is a common verification tool for probabilistic forecasts of meteorological  
<sup>460</sup> (*Gneiting and Raftery*, 2007) and hydrologic (*Laio and Tamea*, 2007; *Thyer et al.*, 2009; *Renard et al.*,  
<sup>461</sup> 2011) variables. This graph helps diagnose errors in ensemble mean (bias) and spread (dispersion) as

462 causes for the deviation from the theoretical 1:1 line for perfect distribution forecasts.  
 463 The uniformity of the  $u$ 's can be tested formally using the Kolmogorov-Smirnov statistic (*Kolmogorov*,  
 464 1933; *Smirnov*, 1948) and we can inspect the  $u$ 's for any sign of non-independence or a trend using the  
 465 uniform condition test (*Cox and Lewis*, 1966). To simplify pairwise comparison of Q-Q plots, we can  
 466 concatenate the deviations of the  $u$ 's from the 1:1 line into a single numerical value or index. For example,  
 467 *Renard et al.* (2010) introduced the so-called reliability  $R_l$  as affine transformation of the taxicab distance  
 468 between the empirical quantiles,  $u_t = F_{P_t}(\omega_t)$ ;  $t = (1, \dots, n)$  and the corresponding quantiles of the  
 469 standard uniform distribution

$$R_l = \frac{2}{n} \sum_{t=1}^n \left| u'_t - \frac{t}{n} \right|, \quad (19)$$

470 where  $u'_1, \dots, u'_n$  denote the ordered values of  $u_1, \dots, u_n$ ,  $|\cdot|$  is the absolute value operator, and the  
 471 true quantiles jump up by  $1/n$  at each of the  $n$  observations. The multiplier of two scales the index  
 472 to the closed interval between 0 (most reliable) and 1 (least reliable). Note that Equation (19) is at  
 473 odds with the formal definition of reliability derived from reliability diagrams of probability forecasts for  
 474 dichotomous events (*Dimitriadis et al.*, 2021). This formal definition is presented in Section 7.

476 The predictive Q-Q plot provides a simple and assumption-free graphical summary of the reliability of  
 477 distribution forecasts. This graph has become commonplace in the hydrologic literature for evaluating  
 478 predictive distributions of precipitation (*Renard et al.*, 2011) and discharge (*Thyer et al.*, 2009; *Renard*  
 479 *et al.*, 2011; *Evin et al.*, 2013), compare, contrast and rank different formulations of the likelihood function  
 480 (*Evin et al.*, 2014; *McInerney et al.*, 2017, 2019) and characterize model input, output and structural  
 481 errors (*Renard et al.*, 2010). But the Q-Q plot and reliability index  $R_l$  should not be used as sole  
 482 determinants of the quality of distribution forecasts (*Gneiting et al.*, 2007; *Renard et al.*, 2011). In a  
 483 thought-provoking example, *Hamill* (2001) demonstrated that the probability integral transform may  
 484 yield a uniform histogram on the unit interval, even if every single forecast is biased. Thus, uniformity of  
 485 the PIT values is a necessary but not a sufficient condition for determining ensemble reliability.

486 To address these limitations, *Gneiting et al.* (2007) proposed a more diagnostic approach to the evaluation  
 487 of predictive performance that is based on maximizing the sharpness of the distribution forecasts subject  
 488 to calibration. Within this context, sharpness refers to the concentration of the predictive distributions  
 489 and is a property of the forecasts only. Calibration refers to the statistical consistency between the  
 490 forecast distributions and the observations and is a joint property of the predictions and the outcomes

491 that materialize. The sharpness principle has a theoretical underpinning under the assumption of  
 492 autocalibration (*Tsyplakov*, 2011) and has become a useful working paradigm for probabilistic forecasting  
 493 and forecast evaluation.

494 Table 1 presents three other measures that have found application and use in hydrology for evaluating  
 the accuracy of probabilistic forecasts. The coefficient of variation  $C_v$  is a dimensionless measure of the

Table 1: Time-averaged performance measures,  $\overline{M}(\mathbf{P}, \boldsymbol{\omega})$ , of distribution forecasts  $\mathbf{P} = \{P_1, \dots, P_n\}$  and verifying observations  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_n)^\top$ .

Performance measure	Symbol	$\overline{M}(\mathbf{P}, \boldsymbol{\omega})$	Reference
Reliability <sup>a</sup>	$R_l$	$\frac{2}{n} \sum_{t=1}^n  u'_t - \frac{t}{n} $	<i>Renard et al.</i> (2010)
Coefficient of variation <sup>b</sup>	$C_v$	$\frac{1}{n} \sum_{t=1}^n \frac{\sigma_{P_t}}{\mu_{P_t}}$	<i>Evin et al.</i> (2013)
Coverage <sup>c,d</sup>	$C$	$\frac{1}{n} \sum_{t=1}^n \mathbb{1}\{l_t \leq \omega_t \leq u_t\}$	<i>Dunsmore</i> (1968)
Width <sup>d</sup>	$W$	$\frac{1}{n} \sum_{t=1}^n (u_t - l_t)$	<i>Raftery et al.</i> (2005)

<sup>a</sup>  $u_t = F_{P_t}(\omega_t)$ ;  $t = (1, \dots, n)$  and  $u$ 's are sorted in ascending order to yield  $u'_1, \dots, u'_n$ ;  
 $F_P(x)$  is the cumulative distribution function of  $P$

<sup>b</sup> Replace with sample mean,  $m_P$ , and sample standard deviation,  $s_P$ , in case of an ensemble forecast

<sup>c</sup> The indicator function  $\mathbb{1}\{a\}$  returns 1 if  $a$  is true and zero otherwise

<sup>d</sup> Lower,  $l_t = F_{P_t}^{-1}(\alpha/2)$ , and upper,  $u_t = F_{P_t}^{-1}(1 - \alpha/2)$ , endpoints of  $100(1 - \alpha)\%$  prediction interval at  $\alpha \in (0, 1)$  significance level;  $F_P^{-1}(x)$  is the inverse cumulative distribution (quantile) function of  $P$

495

496 extent of variability (dispersion) in relation to the mean of the distribution. This measure should only be  
 497 computed for data measured on so-called ratio scales which have a meaningful zero point. This measure  
 498 is related to the conjectured sharpness principle of *Raftery et al.* (2005). Smaller values of the  $C_v$  are  
 499 preferred subject to the intervals having the right coverage. The coverage,  $C$ , equals the fraction of  
 500 observations inside the  $\gamma = 100(1 - \alpha)\%$  prediction intervals. To be statistically meaningful and robust,  
 501  $C$  should equal  $1 - \alpha$  at a significance level  $\alpha \in (0, 1)$ . The width,  $W$ , measures the average size of the  $\gamma\%$   
 502 prediction intervals. Despite their intuitive appeal and ease of interpretation, none of the performance  
 503 metrics of Table 1 provides a complete evaluation of the forecast density, and may even be invariant to the  
 504 *true distribution Q* (see Figure 4). Specifically, the width  $W$  and coefficient of variation  $C_v$  are properties  
 505 of the predictive distribution only and, thus, do not guarantee honest forecasts. This invariance to the  
 506 materialized outcome,  $\omega$ , is easily illustrated using Fig. 4a and b. A shift of the distribution forecasts to

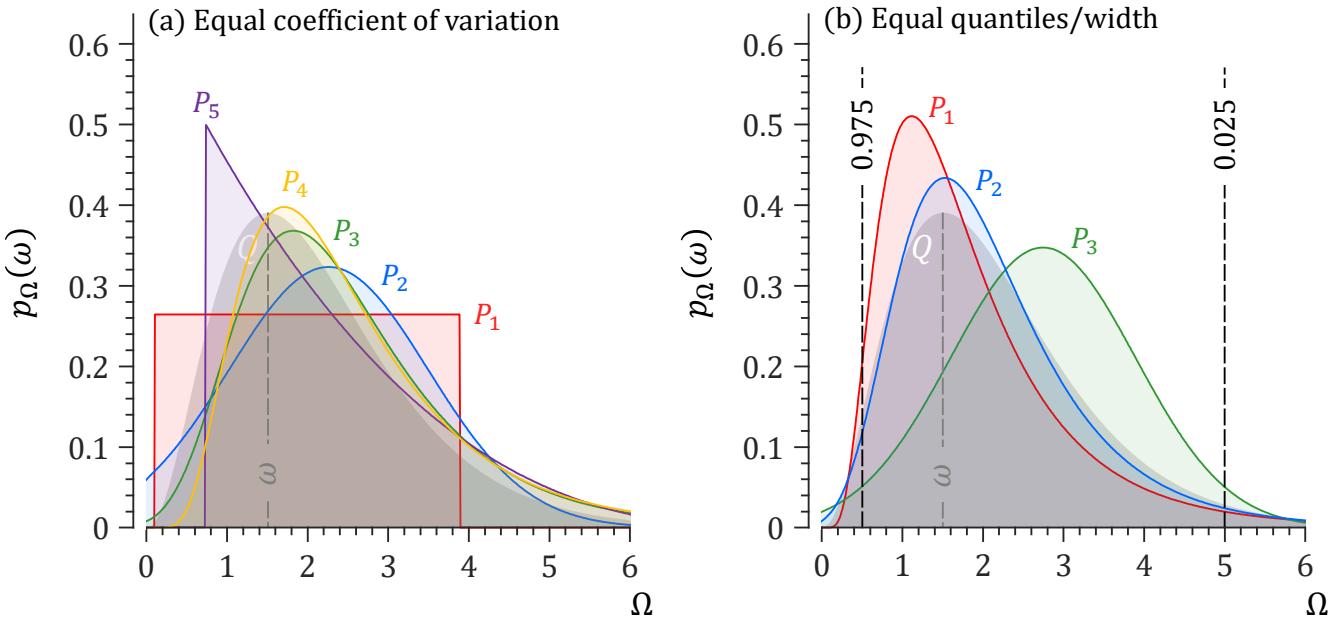


Figure 4: Hypothetical *true distribution*  $Q = \mathcal{G}(3.36, 0.64)$  (gray) and different *probabilistic forecasts*  $P$  (in color) with an equal (a) coefficient of variation,  $C_v = 0.546$ , and (b) width,  $W = 4.5$ , at  $\alpha = 0.05$  using lower and upper quantiles of 0.5 and 5.0, respectively. The peak of the *true distribution* is equal to the verifying observation,  $\omega$ . The distribution forecasts of the left graph are used in later studies:  $P_1 = \mathcal{U}(0.11, 3.89)$ ,  $P_2 = \mathcal{N}(2.26, 1.52)$ ,  $P_3 = \mathcal{GEV}(0.04, 1, 1.86)$ ;  $P_4 = \mathcal{LN}(0.80, 0.26)$  and  $P_5 = \mathcal{GP}(-0.28, 2, 0.73)$ .

507 the right will substantially reduce their overlap with the *true* distribution, but not affect in any way  
 508 their widths and coefficients of variation, which will remain fixed at  $W = 4.5$  and  $C_v = 0.546$ . The  
 509 reliability and coverage measure only two aspects of the statistical consistency between the distributional  
 510 forecasts and the observations. Even if the coverage  $C$  is adequate at a given significance level  $\alpha$ , this  
 511 does not guarantee accurate prediction intervals for other confidence levels (Christoffersen, 1998). This  
 512 necessitates the simultaneous conditional calibration of many different quantile forecasts, which is a  
 513 daunting task.

514 As should be evident from our discussion, the  $R_l$ ,  $C_v$ ,  $C$  and  $W$  performance metrics of Table 1 measure  
 515 different and complementary aspects of the *distribution forecast*  $P$ . This is diagnostically appealing  
 516 (we address this topic in Section 7), but frustrates forecast evaluation as we cannot aggregate the  $R_l$ ,  
 517  $C_v$ ,  $W$  and  $C$  criteria into a single performance index without assigning arbitrary weights. One can  
 518 adopt "Paretoian" theory of general equilibrium and use non-dominated sorting of the performance metrics  
 519 to rank the distribution forecasts (McInerney *et al.*, 2017, 2019). This is a pragmatic solution but  
 520 synonymous to an arbitrary and incomplete evaluation of distribution forecasts.

## 521 5 Scoring Rules

522 Scoring rules are indispensable in our search for the *true* forecast distribution  $Q$ , nevertheless, have not  
 523 yet entered mainstream use in the hydrologic community. In terms of elicitation, the role of scoring rules  
 524 is to encourage the assessor to make careful assessments and to be honest (*Garthwaite et al.*, 2005). In  
 525 terms of evaluation, scoring rules measure the quality of the probabilistic forecasts, reward probability  
 526 assessors for forecasting jobs, and rank competing forecast procedures. Meteorologists refer to this broad  
 527 task as *forecast verification*. We briefly review the most important statistical concepts and explicate our  
 528 notation.

### 529 5.1 Definition

530 A *probabilistic forecast* is any probability measure  $P \in \mathcal{P}$  on  $\Omega$ , the sample space,  $P(\Omega) = 1$ . A *point*  
 531 *forecast* is a functional, i.e. a set-valued mapping  $P \rightarrow T(P) \subseteq \Omega$ , from a class of probability distributions  
 532  $P$  to the real line  $\mathbb{R}$  with the mean or expectation, quantiles and expectiles being key examples. Per  
 533 definition 1 on Page 4, a *scoring function*  $s(x, \omega)$  is any real-valued function  $s : \Omega \times \Omega \rightarrow \mathbb{R}$  which  
 534 quantifies the reward to the forecaster of a point forecast  $x \in \Omega$  and verifying observation  $\omega \in \Omega$ . Then,  
 535 according to definition 2, a *scoring rule*  $S(P, \omega)$  is any extended real-valued and  $\mathcal{P}$ -quasi-integrable  
 536 function  $S : \mathcal{P} \times \Omega \rightarrow \overline{\mathbb{R}} \equiv [-\infty, \infty]$  for all  $P \in \mathcal{P}$  which measures the reward when forecast  $P$  is issued  
 537 and the observation  $\omega \in \Omega$  materializes. The function  $\mathcal{S}(P, Q) : \mathcal{P} \times \mathcal{P} \rightarrow \overline{\mathbb{R}}$  is equal to the expected  
 538 value of  $S(P, \omega)$  under the *true distribution*  $Q$  of  $\omega$ .

### 539 5.2 Theory

540 The expected score of the probabilistic forecast  $P$  under the *true* distribution  $Q \in \mathcal{P}$  is defined as follows

$$541 \quad 542 \quad \mathcal{S}(P, Q) = \mathbb{E}_{\omega \stackrel{\mathcal{D}}{\sim} Q} [S(P, \omega)] = \int_{\Omega} S(P, \omega) dQ(\omega) \quad (20)$$

543 where the symbol  $\stackrel{\mathcal{D}}{\sim}$  means “*distributed according to*”. Note that the order of the arguments has reversed  
 544 with respect to the convention used in information theory. The *probabilistic forecast*  $P$  precedes the *true*  
 545 *distribution*  $Q$ . This is in line with the statistical forecasting literature (*Gneiting and Raftery*, 2007;  
 546 *Bröcker*, 2009). It is our convention that a higher score indicates a good forecast and, thus, our scoring

547 rules  $S(P, Q)$  are positively oriented and defined as reward functions. Then, a scoring rule  $S$  is said to  
 548 be *proper* relative to  $\mathcal{P}$  if  
 549

$$S(P, Q) \leq S(Q, Q) \quad \text{for all } P, Q \in \mathcal{P}, \quad (21)$$

550 and is considered *strictly proper* if the above condition holds with equality if and only if  $P = Q$ . This  
 551 implies that a *strictly proper* score rule is a sufficient condition, whereas a *proper* score rule is a necessary  
 552 but not sufficient condition. In plain words, if  $S(P, Q)$  is a *strictly proper* score rule, then the larger its  
 553 value, the closer the distribution of  $P$  will be to that of  $Q$ . This is not true for *proper* scoring rules,  
 554 which can attain a maximum score even if  $P \neq Q$  (Vrugt *et al.*, 2022). Based on early recommendations  
 555 by Brier (1950) and Shuford *et al.* (1966), we restrict attention to the class of *proper* scoring rules. This  
 556 includes *strictly proper* scoring rules.

557 A scoring rule  $S : \mathcal{P} \times \Omega \rightarrow \overline{\mathbb{R}}$  is *regular* relative to the class  $\mathcal{P}$  if  $S(P, Q)$  is real-valued for all  $P, Q \in \mathcal{P}$ ,  
 558 except possibly that  $S(P, Q) = \infty$  if  $P \neq Q$ . If  $S$  is *regular* and *proper*, then the excess score

$$d(P, Q) = S(Q, Q) - S(P, Q), \quad P, Q \in \mathcal{P}, \quad (22)$$

560 measures the discrepancy of the *probabilistic forecast*  $P \in \mathcal{P}$  from the *true distribution*  $Q \in \mathcal{P}$ . This is  
 561 a *divergence function* alike the relative entropy in Equations (1) and (11) and equal to a measure of the  
 562 difference between two points defined in terms of a continuously-differentiable expected score function,  
 563  $H(P) : \mathcal{P} \rightarrow \mathbb{R}$ . For positively oriented *proper* scoring rules  $H(P)$  is equal to the pointwise supremum  
 564 (least upper bound) over the convex class of probability measures  $Q$  on  $\mathcal{P}$  (Gneiting and Raftery, 2007)

$$\begin{aligned} H(P) &= \sup_{Q \in \mathcal{P}} S(Q, P), \quad P \in \mathcal{P}, \\ &= S(P, P). \end{aligned} \quad (23)$$

568 and is convex on  $\mathcal{P}$  since  $S(Q, P)$  is linear in  $P$  (Rockafellar, 1970). The statement holds with proper  
 569 replaced by strictly proper, and convex replaced by strictly convex. If the sample space is finite and the  
 570 expected score function  $H(P)$  smooth, then  $d(P, Q)$  equals the Bregman (1967) divergence associated  
 571 with the convex function  $H(P)$ .

572 To understand the relationship between  $H(P)$ ,  $S(P, Q)$  and  $d(P, Q)$  Figure 5 displays the entropy function  
 573 of the *strictly proper* quadratic score for a binary event (*rain* or *no rain*)

$$\begin{aligned} H_{\text{QS}}(p) &= p^2 + (1-p)^2 \\ &= 2p^2 - 2p + 1 \\ \implies H_{\text{QS}}(p) &\mapsto p(p-1), \end{aligned} \quad (24)$$

578 where  $p \in [0, 1]$  is the quoted *probability* for *rain* to occur and the symbol  $\mapsto$  signifies affine transformation.  
579 This is equivalent to a sample space  $\Omega = \{1, 0\}$  with probability  $1 - p$  that no ( $= 0$ ) precipitation was  
580 observed. Entropy functions and score divergences of (*strictly*) *proper* scoring rules are invariant to such  
geometric transformation. The expected score function of the quadratic score is strictly convex on  $p$

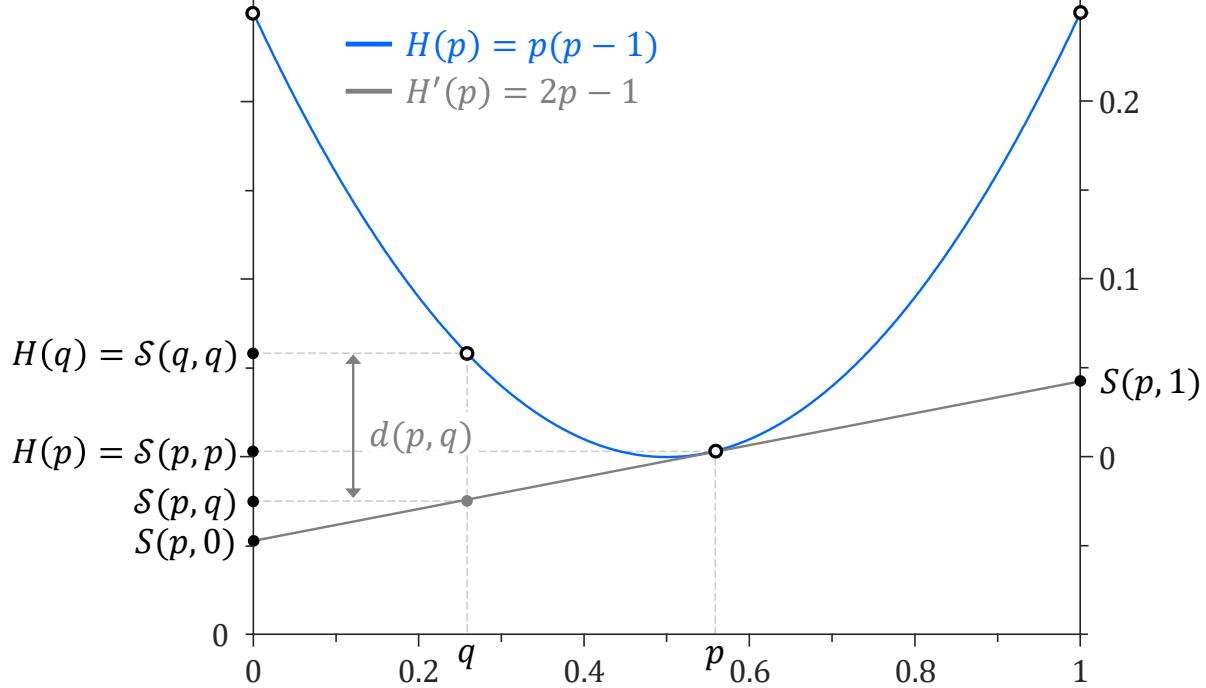


Figure 5: Generalized entropy function  $H(p) = p(p-1)$  (blue curve) of the quadratic score for a dichotomous event  $\Omega = \{1, 0\}$  with *probability forecast*  $(p, 1-p)$  and *true probability*  $(q, 1-q)$  with  $p, q \in [0, 1]$ . We present the values of the quadratic scoring rule  $\mathcal{S}(p, q)$  at  $p$  and  $q$  (solid black dots) and display the so-called Bregman divergence,  $d(p, q)$ . For any probability forecast,  $p \in [0, 1]$ , the expected score,  $\mathcal{S}(p, q) = qS(p, 1) + (1-q)S(p, 0)$ , equals the ordinate of the tangent to  $H$  at  $p$  (solid gray line) when evaluated at  $q \in [0, 1]$ . In particular, the scores,  $S(p, 0) = H(p) - pH'(p)$  and  $S(p, 1) = H(p) + (1-p)H'(p)$ , equal the tangent at  $q = 0$  and  $q = 1$ , respectively. The divergence,  $d_{QS}(p, q) = \mathcal{S}(q, q) - \mathcal{S}(p, q)$ , is equal to the difference between  $H(q)$  and the tangent at  $p$  when evaluated at  $q$  (Adapted after Fig. 1 of Gneiting and Raftery (2007) and Fig. 8 of Buja *et al.* (2005)).

581  
582 and satisfies continuity. The  $H(p)$  function is referred to in the statistical literature as the *information*  
583 *measure* or *(generalized) entropy function* associated with the scoring rule  $S$  (Grünwald and Dawid, 2004;  
584 Buja *et al.*, 2005). This is the maximally achievable utility. Some authors refer instead to  $-H(p)$  as  
585 the entropy function (Bröcker, 2009; Dawid and Musio, 2014) or *coherent uncertainty function* (Dawid  
586 and Sebastiani, 1999). According to Figure 5 and Equation (22), the score divergence  $d(p, q)$  equals the  
587 difference of the value of  $H$  at point  $q$  and the first-order Taylor expansion of  $H$  around point  $p$  evaluated  
588 at point  $q$

$$589 \quad d(p, q) = \mathcal{S}(q, q) - \mathcal{S}(p, q)$$

$$\begin{aligned}
&= H(q) - (H(p) - H'(p)(p - q)) \\
&= H(q) - H(p) + H'(p)(p - q),
\end{aligned} \tag{25}$$

where  $H'(p) = dH(p)/dp$  is the derivative of the entropy function with respect to  $p$ . The squared Euclidean distance  $d(p, q) = \|p - q\|_2^2 = (p - q)^2$  is the canonical example of a Bregman distance, generated by the convex function  $H(p) = p^2$ . Indeed, according to Equation (25) the generalized entropy function  $H(p) = p^2$  has divergence function

$$\begin{aligned}
d(p, q) &= H(q) - H(p) + H'(p)(p - q) \\
&= q^2 - p^2 + 2p(p - q) \\
&= q^2 + p^2 - 2pq \\
&= (p - q)^2
\end{aligned} \tag{26}$$

which is equal to the well-known Brier score (*Brier*, 1950).

### 5.3 Scoring Rules for Categorical Forecasts

For a categorical forecast of a finite number of  $m$  mutually exclusive and collectively exhaustive events  $\Omega = \{1, \dots, m\}$  the *distribution forecast* is a probability vector  $\mathbf{p} = (p_1, \dots, p_m)^\top$  issued on the convex class  $\mathcal{P} = \mathcal{P}_m$  defined in Equation (9). Then Equation (25) may be written as

$$d_{\text{QS}}(\mathbf{p}, \mathbf{q}) = H(\mathbf{q}) - H(\mathbf{p}) + \langle \nabla H(\mathbf{p}), \mathbf{p} - \mathbf{q} \rangle \tag{27}$$

where the gradient  $\nabla H(\mathbf{p}) : \mathbb{R}^m \rightarrow \mathbb{R}^m$  at  $\mathbf{p} \in \mathcal{P}_m$  is a vector-valued function

$$\nabla H(\mathbf{p}) = \frac{\partial H(\mathbf{p})}{\partial \mathbf{p}} = \left( \frac{\partial H(\mathbf{p})}{\partial p_1}, \frac{\partial H(\mathbf{p})}{\partial p_2}, \dots, \frac{\partial H(\mathbf{p})}{\partial p_m} \right)^\top \tag{28}$$

and  $\langle \mathbf{a}, \mathbf{b} \rangle$  denotes the inner product of the  $m$ -vectors  $\mathbf{a}$  and  $\mathbf{b}$ . Furthermore, a *regular* scoring rule of a categorical forecast is *proper* if and only if (*McCarthy*, 1956; *Savage*, 1971)

$$S(\mathbf{p}, j) = H(\mathbf{p}) - \langle \nabla H(\mathbf{p}), \mathbf{p} \rangle + \nabla H_j(\mathbf{p}) \quad \text{for } j = (1, \dots, m) \tag{29}$$

and reduces to a pair of functions,  $S(p, 1) : p \in [0, 1] \rightarrow \overline{\mathbb{R}}$  and  $S(p, 0) : p \in [0, 1] \rightarrow \overline{\mathbb{R}}$ , for the binary event of *rain* or *no rain*. For a probability quote  $p$  the reward of the forecaster will equal to  $S(p, 1)$  if rainfall materializes and  $S(p, 0)$  otherwise. The expected score in Equation (20) then equals  $\mathcal{S}(p, q) = qS(p, 1) + (1 - q)S(p, 0)$ , where  $q$  is the true probability of *rain*. For any two probability assignments  $\mathbf{p}$  and  $\mathbf{q}$  with *true probabilities*  $\mathbf{q} = (q_1, \dots, q_m)^\top$  constrained to the probability simplex,

Table 2: Quadratic, logarithmic and pseudospherical scoring rules for categorical variables: Entropy function, scoring rule, expected score function and divergence function for a *distribution forecast*  $\mathbf{p} = (p_1, \dots, p_m)^\top$  on the convex class  $\mathcal{P} = \mathcal{P}_m$  of  $m \geq 2$  mutually exclusive and collectively exhaustive events,  $\Omega = \{1, \dots, m\}$ . The  $m$ -vector  $\mathbf{q} = (q_1, \dots, q_m)^\top$  lists the true event probabilities.

Name	Score	Entropy function $H(\mathbf{p})$	Scoring rule $S(\mathbf{p}, j)$	Expectation $\mathcal{S}(\mathbf{p}, \mathbf{q})$	Divergence $d(\mathbf{p}, \mathbf{q})$
Quadratic <sup>a,b</sup>		$\sum_{k=1}^m p_k^2$	$2p_j - \sum_{k=1}^m p_k^2$	$2 \sum_{k=1}^m p_k q_k - \sum_{k=1}^m p_k^2$	$\sum_{k=1}^m (p_k - q_k)^2$
Logarithmic <sup>c</sup>		$\sum_{k=1}^m p_k \log_b(p_k)$	$\log_b(p_j)$	$\sum_{k=1}^m q_k \log_b(p_k)$	$\sum_{k=1}^m q_k \log_b\left(\frac{q_k}{p_k}\right)$
Pseudospherical <sup>d,e</sup>		$\ \mathbf{p}\ _\eta^1$	$\frac{p_j^{\eta-1}}{\ \mathbf{p}\ _\eta^{\eta-1}}$	$\frac{\langle \mathbf{p}^{\eta-1}, \mathbf{q} \rangle}{\ \mathbf{p}\ _\eta^{\eta-1}}$	$\ \mathbf{q}\ _\eta^1 - \frac{\langle \mathbf{p}^{\eta-1}, \mathbf{q} \rangle}{\ \mathbf{p}\ _\eta^{\eta-1}}$

<sup>a</sup> Also known as proper linear score. Equals *Brier* (1950) score for a binary event,  $\Omega = \{1, 0\}$

<sup>b</sup> Limiting case ( $\eta \rightarrow 1$ ) of the pseudospherical score if suitably scaled

<sup>c</sup> Remains *strictly proper* under any logarithmic base  $b > 1$ .

<sup>d</sup>  $\|\mathbf{p}\|_b^a = (\sum_{k=1}^m p_k^b)^{a/b}$  and  $\langle \mathbf{p}, \mathbf{q} \rangle = \sum_{k=1}^m p_k q_k$  is the dot product of *forecast* and *true* probabilities

<sup>e</sup> For  $\eta = 2$  the PSS reduces to the spherical score of *Friedman* (1983)

618  $\{\mathbf{q} \in \mathbb{R}_+^{m \times 1} : \mathbf{q}^\top \mathbf{1}_m = 1\}$  and  $\mathbb{R}_+ = [0, \infty)$ , the binary definition of the expected score generalizes to

619 (*Bröcker*, 2009)

$$620 \quad \mathcal{S}(\mathbf{p}, \mathbf{q}) = \sum_{j=1}^m q_j S(\mathbf{p}, j), \quad (30)$$

621 The interpretation of the scoring function  $\mathcal{S}(\mathbf{p}, \mathbf{q})$  is that if  $\omega$  is a random variable of distribution  $Q$ ,

622 then  $\mathcal{S}(\mathbf{p}, \mathbf{q})$  is the mathematical expectation of the score of the assignment  $\mathbf{p}$  in forecasting  $\omega$ .

623 Table 2 presents the entropy functions  $H(\mathbf{p})$  of three commonly used *strictly proper* categorical scoring

624 rules on  $\mathcal{P}_m$ . Next, we use Equations (27), (29) and (30) to derive the mathematical expressions of the

625 scoring rule,  $S(\mathbf{p}, j)$ , expected score function,  $\mathcal{S}(\mathbf{p}, \mathbf{q})$ , and score divergence,  $d(\mathbf{p}, \mathbf{q})$ , of the quadratic,

626 logarithmic and (pseudo)spherical rule presented in Table 2. By definition, the probabilities  $p$  and  $q$  are

627 dimensionless, and, consequently, all functions are unitless, except for the entropy function  $H(\mathbf{p})$  of the

628 logarithmic score which has units of information and, thus, bits if  $b = 2$ .

### 629 5.3.1 Quadratic Score

630 For a categorical forecast of  $m$  events with true probabilities,  $p_1, \dots, p_m$  issued on the probability

631 simplex  $P \in \mathcal{P}_m$  the entropy function of the quadratic score becomes  $H(\mathbf{p}) = \sum_{k=1}^m p_k^2$  and equals an

632 affine transformation of the Gini-index  $G(\mathbf{p}) = \sum_{k=1}^m p_k(1 - p_k)$  (*Gini*, 1909). The gradient  $\nabla H(\mathbf{p}) =$

633  $(2p_1, \dots, 2p_m)^\top$  and Equation (27) leads to the divergence function

$$\begin{aligned}
d_{\text{QS}}(\mathbf{p}, \mathbf{q}) &= H(\mathbf{q}) - H(\mathbf{p}) + \langle \nabla H(\mathbf{p}), \mathbf{p} - \mathbf{q} \rangle \\
&= \sum_{k=1}^m q_k^2 - \sum_{k=1}^m p_k^2 + \langle (2p_1, \dots, 2p_m), (p_1 - q_1, \dots, p_m - q_m) \rangle \\
&= \sum_{k=1}^m q_k^2 - \sum_{k=1}^m p_k^2 + 2\langle \mathbf{p}, \mathbf{p} \rangle - 2\langle \mathbf{p}, \mathbf{q} \rangle \\
&= \sum_{k=1}^m q_k^2 + \sum_{k=1}^m p_k^2 - 2 \sum_{k=1}^m p_k q_k \\
&= \sum_{k=1}^m (p_k - q_k)^2
\end{aligned} \tag{31}$$

640 The scoring rule of the quadratic score is now equal to

$$\begin{aligned}
S_{\text{QS}}(\mathbf{p}, j) &= H(\mathbf{p}) - \langle \nabla H(\mathbf{p}), \mathbf{p} \rangle + \nabla H_j(\mathbf{p}) \\
&= \sum_{k=1}^m p_k^2 - 2\langle \mathbf{p}, \mathbf{p} \rangle + 2p_j \\
&= 2p_j - \sum_{k=1}^m p_k^2
\end{aligned} \tag{32}$$

645 The quadratic score is also known as the proper linear score or *Brier* (1950) scoring rule

$$S_{\text{BS}}(\mathbf{p}, j) = - \sum_{k=1}^m (\delta_{jk} - p_k)^2 = 2p_j - \sum_{k=1}^m p_k^2 - 1, \tag{33}$$

647 where the Kronecker symbol  $\delta_{jk}$  equals one when the  $j^{\text{th}}$  event materializes ( $j = k$ ) and  $\delta_{jk} = 0$  otherwise.

648 With true probabilities,  $\mathbf{q} = (q_1, \dots, q_m)^\top$ , the expected score of the quadratic rule equals

$$\begin{aligned}
\mathcal{S}_{\text{QS}}(\mathbf{p}, \mathbf{q}) &= \sum_{j=1}^m q_j S_{\text{QS}}(\mathbf{p}, j) = \sum_{j=1}^m q_j \left( 2p_j - \sum_{k=1}^m p_k^2 \right) = 2 \sum_{j=1}^m p_j q_j - \sum_{j=1}^m q_j \sum_{k=1}^m p_k^2 \\
&= 2 \sum_{j=1}^m p_j q_j - \sum_{k=1}^m p_k^2
\end{aligned} \tag{34}$$

652 and we can confirm the divergence function of the quadratic score in Table 2 and Equation (31)

$$\begin{aligned}
d_{\text{QS}}(\mathbf{p}, \mathbf{q}) &= \mathcal{S}_{\text{QS}}(\mathbf{q}, \mathbf{q}) - \mathcal{S}_{\text{QS}}(\mathbf{p}, \mathbf{q}) \\
&= 2 \sum_{j=1}^m q_j q_j - \sum_{k=1}^m q_k^2 - 2 \sum_{j=1}^m p_j q_j + \sum_{k=1}^m p_k^2 \\
&= \sum_{k=1}^m q_k^2 + \sum_{k=1}^m p_k^2 - 2 \sum_{k=1}^m p_k q_k \\
&= \sum_{k=1}^m (p_k - q_k)^2.
\end{aligned}$$

658 **5.3.2 Logarithmic Score**

659 The logarithmic score of *Good* (1952) is a linear equivalent of the relative entropy or KL divergence  
 660 (*Lai et al.*, 2011) and also known in the statistical literature as the predictive deviance (*Knorr-Held and*  
 661 *Rainer*, 2001) and ignorance score (*Roulston and Smith*, 2002). The entropy function of the logarithmic  
 662 score  $H(\mathbf{p}) = \sum_{k=1}^m p_k \log_b(p_k)$  is equal to negative Shannon entropy  $-\mathbb{H}(P)$  with gradient function  
 663  $\nabla H(\mathbf{p}) = (\log_b(p_1) + 1, \dots, \log_b(p_m) + 1)^\top$ . The divergence function  $d_{\text{LS}}(\mathbf{p}, \mathbf{q})$  of the logarithmic score  
 664 may now be derived from Equation (27) to yield

$$\begin{aligned} 665 \quad d_{\text{LS}}(\mathbf{p}, \mathbf{q}) &= H(\mathbf{q}) - H(\mathbf{p}) + \langle \nabla H(\mathbf{p}), \mathbf{p} - \mathbf{q} \rangle \\ 666 \quad &= \sum_{k=1}^m q_k \log_b(q_k) - \sum_{k=1}^m p_k \log_b(p_k) + \langle (\log_b(p_1) + 1, \dots, \log_b(p_m) + 1), (p_1 - q_1, \dots, p_m - q_m) \rangle \\ 667 \quad &= \sum_{k=1}^m q_k \log_b(q_k) - \sum_{k=1}^m p_k \log_b(p_k) + \langle \log_b(\mathbf{p}), \mathbf{p} \rangle + \langle \mathbf{1}_m, \mathbf{p} \rangle - \langle \log_b(\mathbf{p}), \mathbf{q} \rangle - \langle \mathbf{1}_m, \mathbf{q} \rangle \\ 668 \quad &= \sum_{k=1}^m q_k \log_b(q_k) - \sum_{k=1}^m q_k \log_b(p_k) + \sum_{k=1}^m p_k - \sum_{k=1}^m q_k \\ 669 \quad &= \sum_{k=1}^m q_k (\log_b(q_k) - \log_b(p_k)) \\ 670 \quad &= \sum_{k=1}^m q_k \log_b\left(\frac{q_k}{p_k}\right), \end{aligned} \tag{35}$$

672 and the scoring rule of the logarithmic score is now equal to

$$\begin{aligned} 673 \quad S_{\text{LS}}(\mathbf{p}, j) &= H(\mathbf{p}) - \langle \nabla H(\mathbf{p}), \mathbf{p} \rangle + \nabla H_j(\mathbf{p}) \\ 674 \quad &= \sum_{k=1}^m p_k \log_b(p_k) - \langle (\log_b(p_1) + 1, \dots, \log_b(p_1) + 1), (p_1, \dots, p_m) \rangle + \log_b(p_j) + 1 \\ 675 \quad &= \sum_{k=1}^m p_k \log_b(p_k) - \sum_{k=1}^m p_k \log_b(p_k) - \sum_{k=1}^m p_k + \log_b(p_j) + 1 \\ 676 \quad &= \log_b(p_j). \end{aligned} \tag{36}$$

678 Thus the logarithmic score has negative Shannon entropy as its generalized entropy function and  
 679 the reverse KL-divergence,  $d_{\text{LS}}(\mathbf{p}, \mathbf{q}) = \sum_{k=1}^m q_k \log_b(q_k/p_k)$ , as its associated score divergence. The  
 680 logarithmic score benefits a strong mathematical-statistical underpinning and links fundamental aspects  
 681 from statistical theory, decision theory and information theory. The order of the arguments of the score  
 682 divergence is reversed with respect to the convention used in information theory, but is in line with  
 683 that of probabilistic forecasting (*Gneiting and Raftery*, 2007; *Bröcker*, 2009). The logarithmic score is  
 684 equal to the logarithmic value of the probability  $p_j$  assigned to the materialized event,  $j$ , and, thus, only

685 strictly proper locally. Roulston and Smith (2002) provide an information-theoretic perspective on the  
 686 logarithmic score and advocate using the so-called ignorance score,  $\mathcal{S}_{\text{IS}}(\mathbf{p}, j) = -\log_b(p_j)$ .

687 With true probabilities,  $\mathbf{q} = (q_1, \dots, q_m)^\top$ , the expected score of the logarithmic rule becomes

$$688 \quad \mathcal{S}_{\text{LS}}(\mathbf{p}, \mathbf{q}) = \sum_{j=1}^m q_j S_{\text{LS}}(\mathbf{p}, j) = \sum_{j=1}^m q_j \log_b(p_j), \quad (37)$$

690 and we can confirm the logarithmic divergence score of Table 2 and Equation (35)

$$\begin{aligned} 691 \quad d_{\text{LS}}(\mathbf{p}, \mathbf{q}) &= \mathcal{S}_{\text{LS}}(\mathbf{q}, \mathbf{q}) - \mathcal{S}_{\text{LS}}(\mathbf{p}, \mathbf{q}) \\ 692 \quad &= \sum_{j=1}^m q_j \log_b(q_j) - \sum_{j=1}^m q_j \log_b(p_j) \\ 693 \quad &= \sum_{k=1}^m q_k \log_b\left(\frac{q_k}{p_k}\right), \\ 694 \end{aligned}$$

695 which again is equal to the relative entropy  $d_{\text{KL}}(\mathbf{q}, \mathbf{p})$  from  $P$  to  $Q$ .

### 696 5.3.3 Pseudospherical Score

697 The entropy function of the pseudospherical score,  $H(\mathbf{p}) = \|\mathbf{p}\|_\eta^1$  is equal to the  $\ell_\eta$ -norm of the forecast  
 698 probabilities  $\mathbf{p} = (p_1, \dots, p_m)^\top$

$$699 \quad \|\mathbf{p}\|_\eta^1 = (p_1^\eta + \dots + p_m^\eta)^{1/\eta}, \quad (38)$$

700 which for  $\eta = 2$  reduces to the Euclidean norm,  $\sqrt{\mathbf{p}^\top \mathbf{p}}$ . The  $m \times 1$ -vector of partial derivatives  $\nabla H(\mathbf{p})$   
 701 now equals

$$702 \quad \nabla H(\mathbf{p}) = \left( \frac{p_1^{\eta-1}}{\|\mathbf{p}\|_\eta^{\eta-1}}, \frac{p_2^{\eta-1}}{\|\mathbf{p}\|_\eta^{\eta-1}}, \dots, \frac{p_m^{\eta-1}}{\|\mathbf{p}\|_\eta^{\eta-1}} \right)^\top, \quad (39)$$

703 where  $\|\mathbf{p}\|_\eta^{\eta-1} = (\sum_{k=1}^m p_k^\eta)^{(\eta-1)/\eta}$ . The gradient vector of the pseudospherical score may thus be written  
 704 as the scalar-vector product,  $\nabla H(\mathbf{p}) = \|\mathbf{p}\|_\eta^{1-\eta} \mathbf{p}^{\eta-1}$ . According to Equation (27) the score divergence  
 705  $d_{\text{PSS}}(\mathbf{p}, \mathbf{q})$  is now equal to

$$\begin{aligned} 706 \quad d_{\text{PSS}}(\mathbf{p}, \mathbf{q}) &= H(\mathbf{q}) - H(\mathbf{p}) + \langle \nabla H(\mathbf{p}), \mathbf{p} - \mathbf{q} \rangle \\ 707 \quad &= \|\mathbf{q}\|_\eta^1 - \|\mathbf{p}\|_\eta^1 + \|\mathbf{p}\|_\eta^{1-\eta} \langle \mathbf{p}^{\eta-1}, \mathbf{p} - \mathbf{q} \rangle \\ 708 \quad &= \|\mathbf{q}\|_\eta^1 - \|\mathbf{p}\|_\eta^1 + \|\mathbf{p}\|_\eta^{1-\eta} \langle \mathbf{p}^{\eta-1}, \mathbf{p} \rangle - \|\mathbf{p}\|_\eta^{1-\eta} \langle \mathbf{p}^{\eta-1}, \mathbf{q} \rangle \\ 709 \quad &= \|\mathbf{q}\|_\eta^1 - \|\mathbf{p}\|_\eta^1 + \|\mathbf{p}\|_\eta^{1-\eta} \|\mathbf{p}\|_\eta^\eta - \|\mathbf{p}\|_\eta^{1-\eta} \langle \mathbf{p}^{\eta-1}, \mathbf{q} \rangle \\ 710 \quad &= \|\mathbf{q}\|_\eta^1 - \|\mathbf{p}\|_\eta^{1-\eta} \langle \mathbf{p}^{\eta-1}, \mathbf{q} \rangle, \end{aligned} \quad (40)$$

712 for  $\mathbf{p}, \mathbf{q} \in \mathcal{P}_m$  and  $\eta > 1$ . The scoring rule of the pseudospherical score now equals

$$\begin{aligned}
713 \quad S_{\text{PSS}}(\mathbf{p}, j) &= H(\mathbf{p}) - \langle \nabla H(\mathbf{p}), \mathbf{p} \rangle + \nabla H_j(\mathbf{p}) \\
714 \quad &= \|\mathbf{p}\|_\eta^1 - \|\mathbf{p}\|_\eta^{1-\eta} \langle \mathbf{p}^{\eta-1}, \mathbf{p} \rangle + \|\mathbf{p}\|_\eta^{1-\eta} p_j^{\eta-1} \\
715 \quad &= \|\mathbf{p}\|_\eta^1 - \|\mathbf{p}\|_\eta^{1-\eta} \|\mathbf{p}\|_\eta^\eta + \|\mathbf{p}\|_\eta^{1-\eta} p_j^{\eta-1} \\
716 \quad &= \|\mathbf{p}\|_\eta^{1-\eta} p_j^{\eta-1}.
\end{aligned} \tag{41}$$

718 With true probabilities,  $\mathbf{q} = (q_1, \dots, q_m)^\top$ , the expected pseudospherical score equals

$$\begin{aligned}
719 \quad \mathcal{S}_{\text{PSS}}(\mathbf{p}, \mathbf{q}) &= \sum_{j=1}^m q_j S_{\text{PSS}}(\mathbf{p}, j) \\
720 \quad &= \sum_{j=1}^m q_j \|\mathbf{p}\|_\eta^{1-\eta} p_j^{\eta-1} \\
721 \quad &= \|\mathbf{p}\|_\eta^{1-\eta} \langle \mathbf{p}^{\eta-1}, \mathbf{q} \rangle,
\end{aligned} \tag{42}$$

723 and we can confirm the divergence function of the pseudospherical score in Table 2 and Equation (40)

$$\begin{aligned}
724 \quad d_{\text{PSS}}(\mathbf{p}, \mathbf{q}) &= \mathcal{S}_{\text{PSS}}(\mathbf{q}, \mathbf{q}) - \mathcal{S}_{\text{PSS}}(\mathbf{p}, \mathbf{q}) \\
725 \quad &= \|\mathbf{q}\|_\eta^{1-\eta} \langle \mathbf{q}^{\eta-1}, \mathbf{q} \rangle - \|\mathbf{p}\|_\eta^{1-\eta} \langle \mathbf{p}^{\eta-1}, \mathbf{q} \rangle \\
726 \quad &= \|\mathbf{q}\|_\eta^1 - \|\mathbf{p}\|_\eta^{1-\eta} \langle \mathbf{p}^{\eta-1}, \mathbf{q} \rangle.
\end{aligned}$$

728 For  $\eta = 2$  the pseudospherical rule in Equation (42) reduces to the well known spherical scoring rule

$$729 \quad \mathcal{S}_{\text{SS}}(\mathbf{p}, \mathbf{q}) = \|\mathbf{p}\|_2^{-1} \langle \mathbf{p}, \mathbf{q} \rangle, \tag{43}$$

731 with associated divergence function

$$732 \quad d_{\text{SS}}(\mathbf{p}, \mathbf{q}) = \|\mathbf{q}\|_2^1 - \|\mathbf{p}\|_2^{-1} \langle \mathbf{p}, \mathbf{q} \rangle. \tag{44}$$

734 We look in more detail at the scoring rules and consider a binary event of *rain* and *no rain* with  
735 *probability forecast*  $(p, 1-p)$  on  $\Omega = \{1, 0\}$ . Table 3 presents mathematical expressions of the scoring  
736 rules for the probability forecast  $p$  of this dichotomous event. When rain materializes,  $j = 1$ , the score is  
737 equal to  $S(p, 1)$ , otherwise  $j = 0$  and the reward is equal to  $S(p, 0)$ . If we now treat the true forecast  
738 probability of the first event,  $q = q_1$ , then the expected score,  $\mathcal{S}(p, q) = qS(p, 1) + (1-q)S(p, 0)$ . Figure  
739 6 displays the expected score of the quadratic (green), logarithmic (red) and spherical (blue) rules as  
740 function of  $p \in [0, 1]$  and  $q = q_1$  (different graphs). We used  $b = 2$  for the logarithmic score to yield  
741  $\mathcal{S}_{\text{LS}}(p, q) = q \log_2(p) + (1-q) \log_2(1-p)$  in units of bits.

Table 3: Strictly proper scoring rules for a dichotomous event (*rain* and *no rain*) with *probability forecast*  $\mathbf{p} = (p, 1 - p)$  on  $\Omega = \{1, 0\}$  with  $p \in [0, 1]$ .

Scoring rule	$S(p, 1)$	$S(p, 0)$
Brier	$-p^2 + 2p - 1$	$-p^2$
Quadratic	$4p - 2p^2 - 1$	$1 - 2p^2$
Logarithmic	$\log_b(p)$	$\log_b(1 - p)$
Spherical	$p(2p^2 - 2p + 1)^{-1/2}$	$(1 - p)(2p^2 - 2p + 1)^{-1/2}$
Pseudospherical	$p^{\eta-1}(p^\eta + (1 - p)^\eta)^{(1-\eta)/\eta}$	$(1 - p)^{\eta-1}(p^\eta + (1 - p)^\eta)^{(1-\eta)/\eta}$

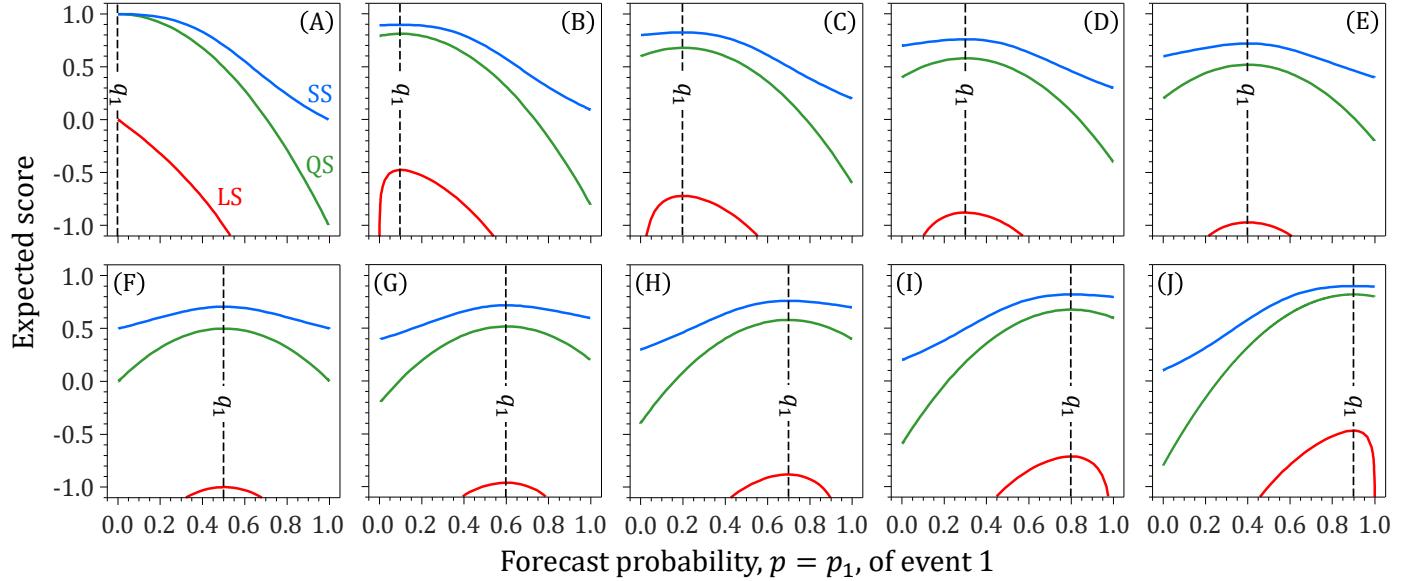


Figure 6: Binary event,  $\Omega = \{1, 0\}$ : Expected value of the quadratic (green), logarithmic (red) and spherical (blue) scoring rules as function of the *true probability*  $q = q_1$  of the first event, (a)  $q = 0$ , (b)  $q = 0.1$ , (c)  $q = 0.2$ , (d)  $q = 0.3$ , (e)  $q = 0.4$ , (f)  $q = 0.5$ , (g)  $q = 0.6$ , (h)  $q = 0.7$ , (i)  $q = 0.8$ , and (j)  $q = 0.9$ .

742 The three scoring rules differ in their response to the quoted forecast probability  $p = p_1$  of the *rain* event.  
 743 The colored lines do not intersect and have a dissimilar functional shape, magnitude and range. But  
 744 despite these differences, the three scoring rules have one property in common. The expected scores of  
 745 QS, LS and SS are always maximized at the true *rain* probability  $q = q_1$ . In other words, the forecaster's  
 746 reward is largest when he/she quotes  $p = q$ . This is exactly what (strict) propriety of the scoring rules  
 747 implies and prompts a forecaster to be honest and report the *true* probabilities. The expected score  
 748 decreases with increasing distance between the quoted and true *rain* probabilities  $p$  and  $q$ , respectively.  
 749 The rate of decline is largest for the LS followed by the QS and SS. Note that magnitude differences are

750 inconsequential as *strictly proper* scoring rules  $S$  remain *strictly proper* under affine transformation

751 
$$S^*(\mathbf{p}, j) = cS(\mathbf{p}, j) + \hbar(j), \quad (45)$$

752 where  $c \neq 0$  is a constant and  $\hbar(j)$  is a  $\mathcal{P}$ -integrable function (*Gneiting and Raftery*, 2007). If  $c < 0$  the  
753 orientation of  $S^*(\mathbf{p}, j)$  changes from a reward to a loss function.

754 **5.4 Numerical examples**

755 While it is generally agreed upon that scoring rules must at least be proper to adequately quantify the  
756 accuracy of probabilistic forecasts (*Winkler et al.*, 1996; *Gneiting and Ranjan*, 2011), the question of  
757 which ones to use in a specific applications remains largely open (*Gneiting and Raftery*, 2007; *Alexander  
et al.*, 2022). For the time being, we restrict our attention to the three categorical scoring rules of Table 2.

759 **5.4.1 Case Study I: Simple Illustration**

760 We revisit the distribution forecasts of Fig. 4a and turn the PDF of the *true* forecast distribution  $Q$  with  
761 continuous sample space  $\Omega = [0, 6]$  into a PMF using  $m = 60$  equally spaced values,  $\omega_i = (6i - 3)/m$ ,  
762 where  $i = (1, \dots, m)$ . The probability of each value (event) is determined from the CDF of  $Q$  and make  
763 up the  $m$ -vector  $\mathbf{q} = (q_1, \dots, q_m)^\top$  of *true* probabilities with unit sum. Similarly, we yield the probability  
764 assignment  $\mathbf{p} = (p_1, \dots, p_m)^\top$  for each distribution forecast,  $P_1, \dots, P_5$ . Table 4 lists the generalized  
765 entropy,  $H(\mathbf{p})$ , expectation,  $\mathcal{S}(\mathbf{p}, \mathbf{q})$  and divergence,  $d(\mathbf{p}, \mathbf{q})$ , of the *strictly proper* quadratic, logarithmic  
766 and spherical scoring rules for  $P_1, \dots, P_5$ .

767 The tabulated values of  $H(\mathbf{p})$ ,  $\mathcal{S}(\mathbf{p}, \mathbf{q})$  and  $d(\mathbf{p}, \mathbf{q})$  vary among the scoring rules and distribution forecasts  
768 and confirm several earlier points, (i) the expected score  $\mathcal{S}(\mathbf{p}, \mathbf{q})$  of each scoring rule is maximized when  
769 the forecaster quotes the *true* probabilities, (ii) the logarithmic score is unbounded and operates on  
770 the extended real-line  $\overline{\mathbb{R}}$  as it applies an indefinitely large penalty to  $P_1$  and  $P_5$  for each realized event  
771 *a priori* thought impossible by their respective uniform and generalized Pareto distribution forecasts,  
772 (iii) the QS, LS and SS divergence scores  $d(\mathbf{p}, \mathbf{q})$  are strictly positive and zero only when  $P = Q$ , and  
773 (iv) *strictly proper* scoring rules do not necessarily yield the same ranking of the distribution forecasts.  
774 This justifies the use of multiple *strictly proper* scoring rules *Vrugt et al.* (2022). As a reminder, we  
775 reiterate that the logarithmic score (i) has negative Shannon entropy  $-\mathbb{H}(\mathbf{p})$  in Equation (10) as its

Table 4: Entropy,  $H(\mathbf{p})$ , expectation,  $\mathcal{S}(\mathbf{p}, \mathbf{q})$  and divergence,  $d(\mathbf{p}, \mathbf{q})$  of the *strictly proper* categorical scoring rules of Table 2 for distribution forecasts  $P_1, \dots, P_5$  displayed in Fig. 4a using  $m = 60$  discrete values,  $\Omega = \frac{1}{20}\{1, 3, 5, \dots, 117, 119\}$ . Column  $R$  lists the rank of the distribution forecasts. The bottom row presents the values for a perfect distribution forecast,  $P = Q$ .

Forecast	Quadratic Score				Logarithmic Score, $\mathfrak{b} = 2$				Spherical Score			
	$H(\mathbf{p})$	$\mathcal{S}(\mathbf{p}, \mathbf{q})$	$d(\mathbf{p}, \mathbf{q})$	$R^{\$}$	$H(\mathbf{p})$	$\mathcal{S}(\mathbf{p}, \mathbf{q})$	$d(\mathbf{p}, \mathbf{q})$	$R$	$H(\mathbf{p})$	$\mathcal{S}(\mathbf{p}, \mathbf{q})$	$d(\mathbf{p}, \mathbf{q})$	$R$
	Eq.	(34)	(31)			(37)	(35)			(43)	(44)	
$P_1$	0.026	0.022	0.005	5	-5.248	$-\infty$	$\infty$	4	0.162	0.150	0.016	5
$P_2$	0.024	0.025	0.003	4	-5.522	-5.493	0.102	2	0.156	0.158	0.009	4
$P_3$	0.026	0.027	0.001	1	-5.483	-5.438	0.046	1	0.161	0.163	0.003	1
$P_4$	0.028	0.026	0.001	2	-5.389	-5.542	0.150	3	0.167	0.162	0.004	2
$P_5$	0.029	0.025	0.003	3	-5.314	$-\infty$	$\infty$	5	0.172	0.159	0.007	3
$Q$	0.028	0.028	0.000		-5.392	-5.392	0.000		0.166	0.166	0.000	

$\$$  Rank of each distribution forecast obtained from sorting  $d(\mathbf{p}, \mathbf{q})$  in ascending order

776 entropy function and (ii) relative entropy  $d_{\text{KL}}(\mathbf{q}, \mathbf{p})$  in Equation (11) as its divergence score  $d_{\text{LS}}(\mathbf{p}, \mathbf{q})$   
777 but is sometimes criticized for its unboundedness.

778 The QS, LS and SS may not give the exact same ranking of the distribution forecasts, but they are  
779 unanimous in their assessment of  $P_3$  as the best forecast of the *true* distribution  $Q$ . This conclusion is  
780 supported by visual inspection of the distribution forecasts with the lognormal distribution forecast  $P_4$   
781 (yellow) as a close second contender. The results further demonstrate that (i) the entropy  $H(\mathbf{p})$  cannot  
782 be used as sole determinant of the accuracy of a forecast distribution. This is implicit as the entropy is a  
783 function of the forecast distribution only, and (ii) outcomes with a zero-probability do not count in the  
784 computation of the entropy of the logarithmic score in accordance with the limit,  $\lim_{x \downarrow 0} x \log_{\mathfrak{b}}(x) = 0$  for  
785  $\mathfrak{b} > 0$ . This explains the elevated values of  $H(\mathbf{p})$  for  $P_1$  and  $P_5$  under the logarithmic score.

#### 786 5.4.2 Case Study II: Rainfall data

787 Next, we illustrate the application of the scoring rules to 24-h forecasts of daily precipitation probability  
788 from the Finnish Meteorological Institute for the city of Tampere in south-central Finland. The original  
789 data set of  $n = 346$  daily rainfall forecasts is presented in Appendix E. *Hughes and Topp* (2015) used the  
790 individual daily forecasts to provide a diagrammatic interpretation of the Brier scoring rule in Equation  
791 (33) and associated score divergences. We simplify this analysis and report in Table 5 the entropy,  
792 expected score and score divergence of the quadratic, logarithmic and spherical scoring rules for the  
793 normalized probabilities of *true* and *forecasted* rainfall listed in the right block of Table E.1. The bottom

row presents the function values when the forecaster quotes the *true* rainfall probabilities,  $\mathbf{p} = \mathbf{q}$ . The

Table 5: Entropy,  $H(\mathbf{p})$ , expectation,  $\mathcal{S}(\mathbf{p}, \mathbf{q})$  and divergence,  $d(\mathbf{p}, \mathbf{q})$  of the *strictly proper* categorical scoring rules of Table 2 for the *true* and *forecasted* rainfall probabilities of Table E.1.

Forecast	Quadratic Score			Logarithmic Score, $b = 2$			Spherical Score		
	$H(\mathbf{p})$	$\mathcal{S}(\mathbf{p}, \mathbf{q})$	$d(\mathbf{p}, \mathbf{q})$	$H(\mathbf{p})$	$\mathcal{S}(\mathbf{p}, \mathbf{q})$	$d(\mathbf{p}, \mathbf{q})$	$H(\mathbf{p})$	$\mathcal{S}(\mathbf{p}, \mathbf{q})$	$d(\mathbf{p}, \mathbf{q})$
	Eq.	(34)	(31)		(37)	(35)		(43)	(44)
$\mathbf{p}$	0.1241	0.1034	0.0043	-3.1557	-3.3505	0.0451	0.3523	0.3229	0.0052
$\mathbf{q}$	0.1077	0.1077	0.0000	-3.3054	-3.3054	0.0000	0.3282	0.3282	0.0000

794

795 tabulated data confirm our earlier points and are mainly listed for benchmark purposes.

## 796 6 Scoring rules for density forecasts

797 The task of determining whether the *probabilistic forecast*  $P$  matches the *true distribution*  $Q$  appears  
 798 difficult, perhaps hopeless, because  $Q$  is never observed, even after the fact. But as Diebold *et al.* (1998)  
 799 realized early on, the challenges posed by these subtleties are not insurmountable. The scoring rules for  
 800 categorical variables can be generalized to density forecasts to assist in forecast verification of continuous  
 801 variables.

802 Scoring rules for density forecasts are defined up to a so-called Lebesgue measure  $\mu$ . The Lebesgue  
 803 integral plays an important role in probability theory but this topic is hardly taught in mathematics  
 804 courses. To explain the Lebesgue measure, please consider Figure 7 which displays an example Lebesgue  
 805 density of the standard normal distribution. The Riemann integral partitions the domain of a function  
 806 into a collection of small intervals and bars are constructed to meet the height of the graph. Then in  $\mathbb{R}^2$   
 807 the resulting rectangles make up the area under the graph. The Lebesgue integral also uses rectangles,  
 808 but these rectangles are formed by partitioning the function's range (also called codomain) into different  
 809 intervals. For each horizontal slice, a rectangle is drawn with height of the corresponding function  
 810 value and width equal to the length of all intervals on the real line  $\mathbb{R}$  (e.g. sample space) where the  
 811 function reaches approximately this height. This horizontal slicing of the codomain leads to much more  
 812 complicated sets of  $\omega$  values, certainly for multivariate densities. Thus, the Lebesgue definition extends  
 813 integral calculation to a much broader class of functions. Now, the Lebesgue measure  $\mu$  is equal to the  
 814 width of each slice, which, in turn, is the sum of the widths of all rectangles with the same height. For

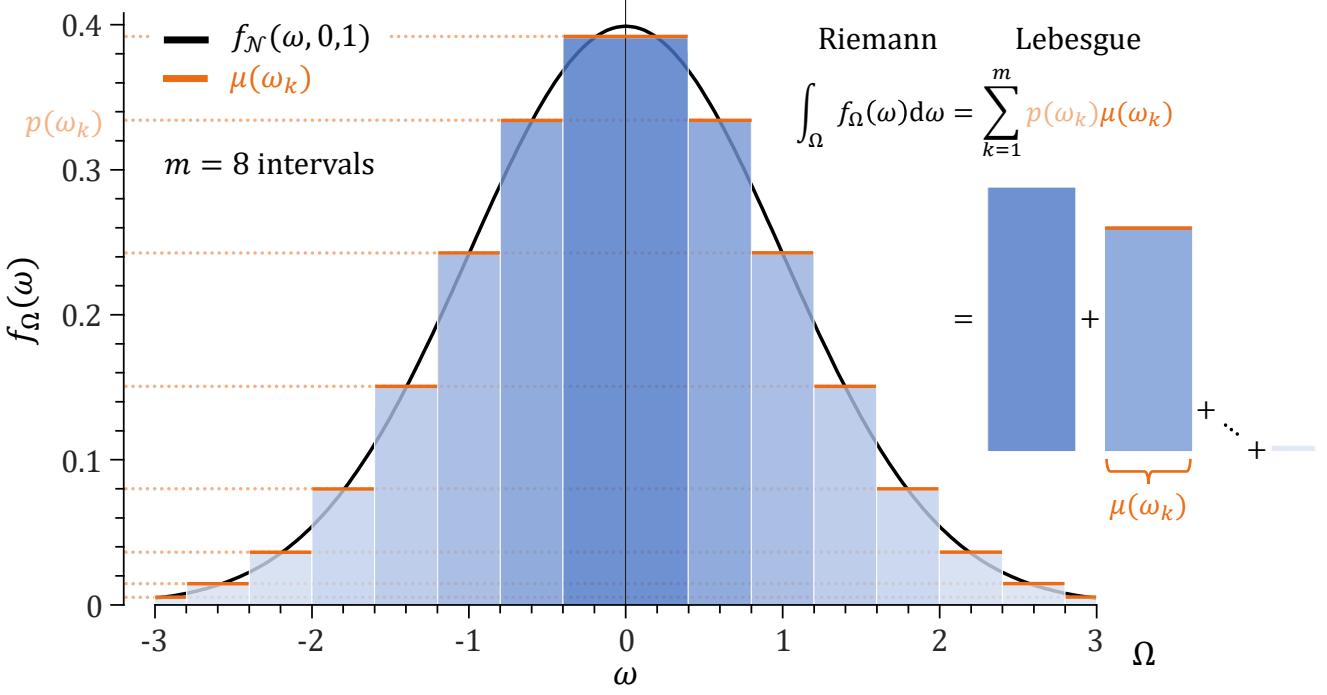


Figure 7: Illustration of the standard normal Lebesgue density on sample space  $\Omega = [-3, 3]$  with range (codomain) of  $N(0, 1)$  partitioned into  $m = 8$  small intervals. We use color coding to discern the intervals of  $\omega$ . The Lebesgue measure  $\mu(\omega_k)$  is equal to the length of each color coded interval containing  $\omega_k$ . The Lebesgue density  $f(\omega)$  is constant in each interval  $\omega_1, \dots, \omega_m$  of  $\omega$  values. The sum of the areas of the rectangles is equal to the Lebesgue integral.

815 univariate densities, we can divide each interval of  $\omega$  in non-overlapping bins with Lebesgue measure the  
 816 bin width. This representation of the Lebesgue density is almost equal to a PMF with each bin made up  
 817 of different points (univariate case) or sets (bivariate case and higher).

## 818 6.1 Univariate forecasts

819 We follow the formal measure-theoretic definition of *Gneiting and Raftery* (2007). Let  $\mu$  be a nonnegative,  
 820 countably additive set function on the measurable space  $(\Omega, \Sigma)$  and let  $\mathcal{L}_\eta(\Omega)$  with  $\eta \in [1, \infty)$  denote the  
 821 class of probability measures  $f : \Omega \rightarrow \mathbb{R}$  on  $(\Omega, \Sigma)$  that are absolutely continuous with respect to the  
 822 measure  $\mu$  on  $\Sigma$  and have an integral

$$823 \|f\|_\eta \equiv \left( \int_{\Omega} f(\omega)^\eta \mu(d\omega) \right)^{1/\eta} < \infty, \quad (46)$$

824 equal to the  $\ell_\eta$ -norm of the density  $f$ . The  $\mu$ -density  $f_P$  of the probabilistic forecast  $P \in \mathcal{L}_\eta$  is  
 825 called a predictive density or density forecast. The above norm is invariant to changes in the *true*  
 826 *distribution Q* that leave the probability of  $\omega$  unchanged and induces a nonnegative metric (divergence)

827  $d(P, Q) = \|f_P - f_Q\|_\eta$  which is zero only if  $P = Q$ .

828 In most applications, the probabilistic forecast  $P$  will consist of a  $m$ -member ensemble  $\mathbf{y} = (y_1, \dots, y_m)^\top$   
 829 at given space and time coordinates. The approximate density of  $\omega$  in a  $\varepsilon$ -neighborhood of a point  $y$  in  
 830  $\mathbb{R}^m$  as

$$831 \quad f(y) = \lim_{\varepsilon \rightarrow 0} \frac{\mu(\omega \cap B_\varepsilon(y))}{\mu(B_\varepsilon(y))} \quad (47)$$

832 where  $B_\varepsilon(y) = \{y \in \mathbb{R} : |\omega - y| < \varepsilon\}$  is the closed ball (sphere) of radius  $\varepsilon$  centered at point  
 833  $y$ . If possible, we compute the unique density of  $\omega$  according to *Bernardo* (1979) (P. 689)  $f(\omega) =$   
 834  $\lim_{\varepsilon \rightarrow 0} P(B_\varepsilon(\omega))/\mu(B_\varepsilon(\omega))$ . The norm  $\|f\|_\eta$  can be expressed as a sequence and/or vector by using the  
 835 counting measure

$$836 \quad \|f_P\|_\eta = \left[ \sum_{k=1}^m \left( \frac{f_P(x_k)}{\sum_{j=1}^m f_P(x_j)} \right)^\eta \right]^{1/\eta}, \quad (48)$$

837 where  $(x_1, \dots, x_m)$  are  $m$  events for which we must evaluate the predictive density. The reciprocal of  
 838 the denominator equals the Lebesgue measure (e.g. histogram bin width) of each event  $x_k$  and assumes  
 839 that the  $x$ 's are evenly distributed on the real line,  $\mathbb{R}$ . This is a pragmatic definition as we estimate  
 840 the predictive density  $f_P$  of the probabilistic forecast  $P$  via kernel smoothing. This method returns  
 841 a  $m$ -vector of equally-spaced  $x$  values along with estimates of their probability density  $f_P(x)$  so that  
 842  $\int_{\mathbb{R}} f_P(x) dx = 1$ .

#### 843 6.1.1 Quadratic, Logarithmic and (Pseudo)spherical scoring rules

844 Scoring rules for the density forecast  $f$  assign a numerical score based on the predictive distribution  $P$   
 845 and on the event or value  $\omega$  that materializes. In analogy to Equation (32), the quadratic or *Brier* (1950)  
 846 score becomes

$$847 \quad S_{\text{QS}}(P, \omega) = 2f_P(\omega) - \|f_P\|_2^2, \quad (49)$$

848 where  $\|f_P\|_2^2$  equals the sum of the squared normalized densities of the forecast distribution  $P$ . This  
 849 scoring is *strictly proper* relative to the class  $\mathcal{L}_2$  and has entropy function,  $H(P) = \|f_P\|_2^2$  and divergence  
 850 function,  $d_{\text{QS}}(P, Q) = \|f_P - f_Q\|^2$ , where  $f_Q$  signifies the density of the true distribution  $Q$ . The  
 851 logarithmic score of Equation (36) equals

$$852 \quad S_{\text{LS}}(P, \omega) = \log_b(f_P(\omega)), \quad (50)$$

853 and is *strictly proper* relative to the class  $\mathcal{L}_1$  of probability measures with entropy function the negative  
 854 Shannon entropy  $H(P) = -\mathbb{H}(P)$  and divergence score the reverse KL-divergence. The logarithmic score

855 is a limiting case of the pseudospherical score for  $\eta \rightarrow 1$

$$856 \quad S_{\text{SPSS}}(P, \omega) = \frac{f_P(\omega)^{\eta-1}}{\|f_P\|_\eta^{\eta-1}}, \quad (51)$$

857 which is *strictly proper* relative to the class  $\mathcal{L}_\eta$ . The strict convexity of the associated entropy function,  
858  $H(P) = \|f_P\|_\eta$ , and the nonnegativity of the divergence function are straightforward consequences of the  
859 Hölder and Minkowski inequalities. For  $\eta = 2$ , we yield the spherical score

$$860 \quad S_{\text{SS}}(P, \omega) = \frac{f_P(\omega)}{\|f_P\|_2}, \quad (52)$$

861 is *strictly proper* relative to the class  $\mathcal{L}_2$  of probability measures.

### 862 6.1.2 Linear scoring rule: impropriety

863 To demonstrate the importance of scoring rule propriety we briefly discuss the intuitively appealing but  
864 improper linear score

$$865 \quad S_{\text{LinS}}(P, \omega) = f_P(\omega), \quad (53)$$

866 Let  $f_P = \frac{1}{2\epsilon}$  and  $f_Q = \frac{1}{\sqrt{2\pi}} \exp(-\omega^2/2)$  denote the Lebesgue densities of the uniform *forecast distribution*  
867  $P$  and standard normal *true distribution*  $Q$  on the interval  $[-\epsilon, \epsilon]$ . According to Equation (20), the  
868 expected score of the probabilistic forecasts  $P$  or  $Q$  under the true distribution  $Q$  becomes

$$\begin{aligned} 869 \quad \mathcal{S}_{\text{LinS}}(P, Q) &= \int_{-\epsilon}^{\epsilon} S_{\text{LinS}}(P, \omega) dQ(\omega) & \mathcal{S}_{\text{LinS}}(Q, Q) &= \int_{-\epsilon}^{\epsilon} S_{\text{LinS}}(Q, \omega) dQ(\omega) \\ 870 &= \int_{-\epsilon}^{\epsilon} \frac{1}{2\epsilon} \frac{1}{\sqrt{2\pi}} \exp(-\omega^2/2) d\omega & &= \int_{-\epsilon}^{\epsilon} \left( \frac{1}{\sqrt{2\pi}} \exp(-\omega^2/2) \right)^2 d\omega \\ 871 &= \frac{1}{2\epsilon} \frac{1}{\sqrt{2\pi}} \int_{-\epsilon}^{\epsilon} \exp(-\omega^2/2) d\omega & &= \frac{1}{2\pi} \int_{-\epsilon}^{\epsilon} \exp(-\omega^2) d\omega \\ 872 &= \frac{1}{2\epsilon} \frac{1}{\sqrt{2\pi}} \left| \sqrt{2\pi} \operatorname{erf}(\omega/\sqrt{2}) \right|_0^\epsilon & &= \frac{1}{2\pi} \left| \sqrt{\pi} \operatorname{erf}(\omega) \right|_0^\epsilon \\ 873 &= \frac{1}{2\epsilon} \operatorname{erf}(\epsilon/\sqrt{2}) & &= \frac{1}{2\sqrt{\pi}} \operatorname{erf}(\epsilon), \end{aligned} \quad (54)$$

875 and the score divergence  $d_{\text{LinS}}(P, Q)$  is equal to

$$\begin{aligned} 876 \quad d_{\text{LinS}}(P, Q) &= \mathcal{S}_{\text{LinS}}(Q, Q) - \mathcal{S}_{\text{LinS}}(P, Q) \\ 877 &= \frac{1}{2\sqrt{\pi}} \operatorname{erf}(\epsilon) - \frac{1}{2\epsilon} \operatorname{erf}(\epsilon/\sqrt{2}). \end{aligned} \quad (55)$$

879 For an infinitesimal interval  $\epsilon \rightarrow 0$  we yield  $d_{\text{LinS}}(P, Q) = -(2\pi)^{-1/2}$ . Figure 8 displays the divergence  
880 of the linear score as function of the Lebesgue measure  $0 < \epsilon \leq 3$ . The score divergence is negative for  
881 small values of  $\epsilon$  and changes sign at the root,  $\epsilon = 1.6221$ , of Equation (55). This demonstrates that  
882 the score divergence of the linear score  $S_{\text{LinS}}(P, \omega) = f_P(\omega)$  does not have a proper zero point. Thus,  
883  $S_{\text{LinS}}(P, \omega)$  is an improper scoring rule.

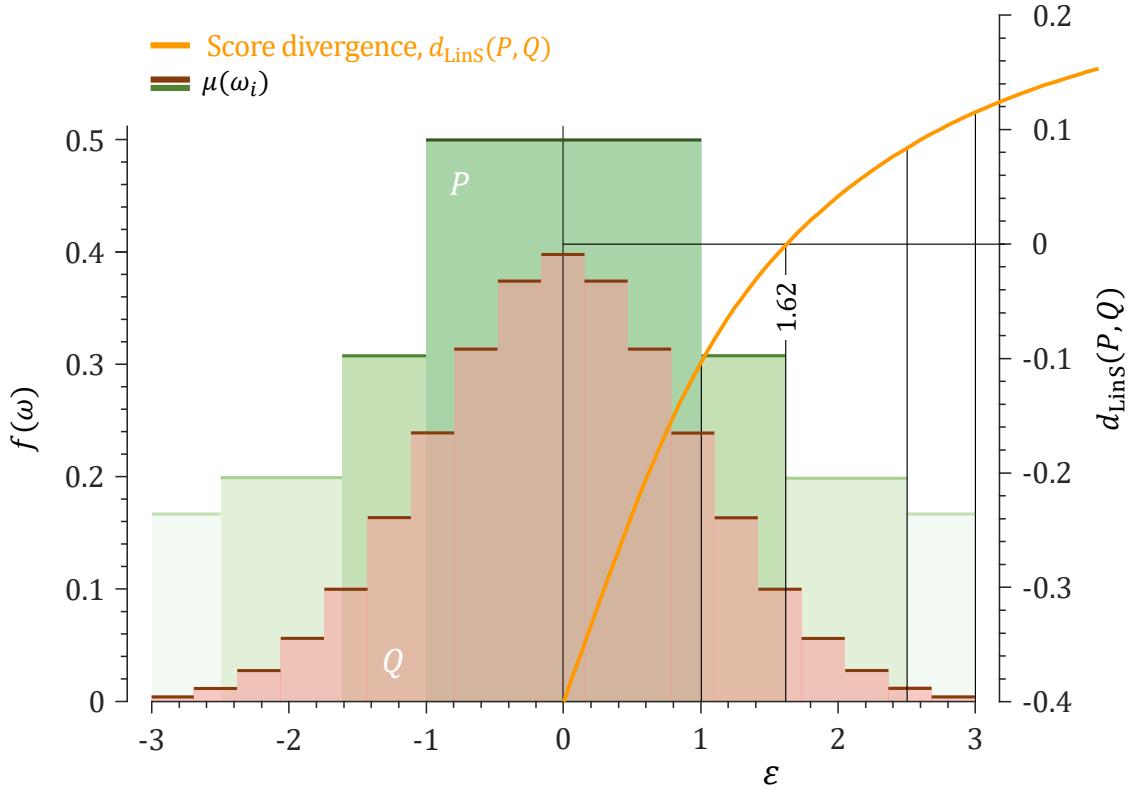


Figure 8: Score divergence  $d_{\text{LinS}}(P, Q)$  of the linear scoring rule  $S_{\text{LinS}}(P, \omega) = f_P(\omega)$  for a uniform probabilistic forecast  $P$  under a standard Gaussian true distribution  $Q$  and symmetric interval  $[-\epsilon, \epsilon]$  with  $\epsilon \in (0, 3]$ . The horizontal green and blue lines correspond to the Lebesgue measure  $\mu(\omega_i)$  or the length of the interval containing event  $\omega_i$ .

### 6.1.3 Continuous Ranked Probability Score

The aforementioned scores are not particularly sensitive to distance in that they do not receive credit when assigning high probabilities to values near but not equal to the materialized outcome. The continuous ranked probability score or CRPS of Figure 9 addresses this deficiency. This scoring rule has found widespread application in the atmospheric sciences and is equal to the integral of the squared distance between the CDF,  $F_P$ , of the distribution forecast  $P$  and the empirical CDF of the observation (*Matheson and Winkler, 1976; Hersbach, 2000*)

$$\begin{aligned}
 S_{\text{CRPS}}(P, \omega) &= - \int_{-\infty}^{\infty} (F_P(z) - \mathbb{1}\{\omega \leq z\})^2 dz \\
 &= - \int_{-\infty}^{\omega} F_P^2(z) dz - \int_{\omega}^{\infty} (F_P(z) - 1)^2 dz,
 \end{aligned} \tag{56}$$

where the indicator function  $\mathbb{1}\{a\}$  returns 1 if  $a$  is true and zero otherwise, and the minus sign reverses the orientation to a reward function. The CRPS is *strictly proper* relative to the subclass  $\mathcal{P}_1 \in \mathcal{P}$  of Borel probability measures that have finite first moment. The CRPS may also be written as integral of

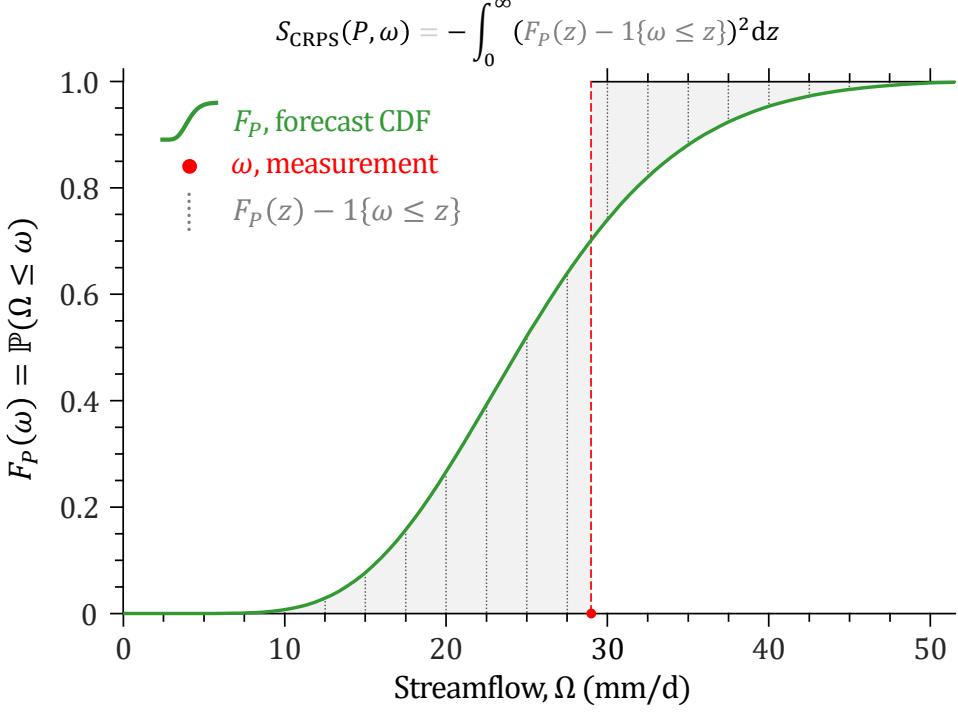


Figure 9: Graphical explanation of the continuous ranked probability scoring rule for a hypothetical streamflow forecast CDF  $F_P$  (black line) and verifying discharge measurement  $\omega$  (red dot). The CRPS is equal to the integral of the squared differences (= gray dotted lines) of  $F_P$  and the Heaviside step function,  $\mathbb{1}\{\omega \leq z\}$ .

897 the Brier probability score of the induced distribution forecast  $F_P(z) = \int_{-\infty}^z f_P(t) dt$  of the binary event  
 898  $\{\omega \leq z\}$  over all thresholds  $z \in \mathbb{R}$

$$899 \quad S_{\text{CRPS}}(P, \omega) = \int_{-\infty}^{\omega} S_{\text{BS}}(F_P(z), \mathbb{1}\{\omega \leq z\}) dz, \quad (57)$$

900 where

$$901 \quad S_{\text{BS}}(F_P(z), \mathbb{1}\{\omega \leq z\}) = -(F_P(z) - \mathbb{1}\{\omega \leq z\})^2. \quad (58)$$

902 We can also write the CRPS using the  $\tau \in [0, 1]$ -quantile forecast  $y_\tau = F_P^{-1}(\tau)$  of  $P$  as follows (*Laio and*  
 903 *Tamea, 2007; Gneiting and Ranjan, 2011; Grushka-Cockayne et al., 2017*)

$$904 \quad S_{\text{CRPS}}(P, \omega) = -2 \int_0^1 (\mathbb{1}\{\omega < y_\tau\} - \tau)(y_\tau - \omega) d\tau, \quad (59)$$

906 with integrand the piecewise linear quantile score (*Friederichs and Hense, 2007; Bracher et al., 2021*)

$$907 \quad S_{\text{QNT}}^\tau(P, \omega) = -2(\mathbb{1}\{\omega < y_\tau\} - \tau)(y_\tau - \omega), \quad (60)$$

909 which is also known as the pinball-loss, tick-loss or check-loss function. The equivalence of Equations  
 910 (56) and (59) can be established via a change of variables (*Laio and Tamea, 2007*). A more friendly

911 formulation of Equation (59) is (see Appendix F)

$$S_{\text{CRPS}}(P, \omega) = \omega(1 - 2F_P(\omega)) + 2 \int_0^1 \tau F_P^{-1}(\tau) d\tau - 2 \int_{F_P(\omega)}^1 F_P^{-1}(\tau) d\tau, \quad (61)$$

912 and promotes closed-form solutions of the CRPS for parametric distribution forecasts  $P$  with analytic  
913 expressions of their CDF,  $F_P$ , and quantile function,  $F_P^{-1}$  (*Jordan*, 2016; *Villez*, 2017). Appendix G  
914 presents such a derivation for a normal distribution forecast  $P = \mathcal{N}(\mu_P, \sigma_P^2)$  to yield

$$S_{\text{CRPS}}(\mathcal{N}(\mu_P, \sigma_P^2), \omega) = \frac{\sigma_P}{\sqrt{\pi}} - 2\sigma_P^2 f_{\mathcal{N}}(\omega, \mu_P, \sigma_P^2) - (\omega - \mu_P)(2F_{\mathcal{N}}(\omega, \mu_P, \sigma_P^2) - 1), \quad (62)$$

915 where  $f_{\mathcal{N}}(x, \mu, \sigma^2)$  and  $F_{\mathcal{N}}(x, \mu, \sigma^2)$  are the normal PDF and CDF, respectively.

916 Nonparametric distribution forecasts do not admit a closed-form expression for the CRPS and, thus, we  
917 must numerically solve for the integral of  $S_{\text{CRPS}}(P, \omega)$  using quadrature rules (*Stael von Holstein*, 1970;  
918 *Unger*, 1985). We can also resort to Lemma 2.2 of *Baringhaus and Franz* (2004) and use the convenient  
919 kernel representation of the CRPS

$$S_{\text{CRPS}}(P, \omega) = \frac{1}{2} \mathbb{E}_P [|y - y^*|] - \mathbb{E}_P [|y - \omega|], \quad (63)$$

920 where  $y$  and  $y^*$  are samples of the forecast distribution  $P$ . For an ensemble forecast of  $m$  members,  
921  $\mathbf{y} = (y_1, \dots, y_m)^\top$ , Equation (63) simplifies to (*Grimit et al.*, 2006)

$$S_{\text{CRPS}}(P, \omega) = \frac{1}{2m^2} \sum_{i=1}^m \sum_{j=1}^m |y_i - y_j| - \frac{1}{m} \sum_{i=1}^m |y_i - \omega|. \quad (64)$$

922 For a point forecast ( $m = 1$ ) the first term equals zero and the CRPS reduces to the negative absolute error,  
923  $s_{\text{NAE}}(y, \omega) = -|y - \omega|$ . Thus, the CRPS is an extension of the absolute error to a distribution forecast.  
924 The computational complexity  $\mathcal{O}(m^2)$  of Equation (64) can be reduced to a total of  $\mathcal{O}(m \log_b(m))$   
925 operations if the CRPS is evaluated in terms of the quantile loss function of Equation (59) using the  
926 sorted ensemble members (*Murphy*, 1970; *Hersbach*, 2000; *Laio and Tamea*, 2007)

$$S_{\text{CRPS}}(P, \omega) = -\frac{2}{m^2} \sum_{i=1}^m |y_i - \omega| \left( m \mathbb{1}\{\omega \leq y_i\} - i + \frac{1}{2} \right). \quad (65)$$

927 The entropy function or information measure of the CRPS

$$H(P) = - \int_{-\infty}^{\infty} F_P(z) (1 - F_P(z)) dz = -\frac{1}{2} \mathbb{E}_P [|y - y^*|], \quad (66)$$

928 coincides with the negative selectivity function of *Matheron* (1984) and the *Gini* (1909) index. The  
929 CRPS divergence function

$$d_{\text{CRPS}}(P, Q) = \int_{-\infty}^{\infty} (F_P(z) - F_Q(z))^2 dz, \quad (67)$$

930 is symmetric by virtue of the quadratic term.

Table 6: Summary of *strictly proper* scoring rules for a density forecast  $f_P$  and verifying observation  $\omega$ . The numerical form assumes that the forecast distribution  $P$  is a  $m$ -member ensemble  $(y_1, \dots, y_m)^\top$ .

Score Name	XX	Scoring rule, $S_{XX}(P, \omega)$		Note
		Analytic	Numerical	
Quadratic	QS	$2f_P(\omega) - \int_{-\infty}^{\infty} f_P^2(y) dy$	$2f_P(\omega) - \sum_{k=1}^m \left( \frac{f_P(y_k)}{\sum_{j=1}^m f_P(y_j)} \right)^2$	a
Logarithmic	LS	$\log_b(f_P(\omega))$	$\log_b(f_P(\omega))$	a
Cnt. Rnk. Prb.	CRPS	$-\int_{-\infty}^{\infty} (F_P(z) - \mathbb{1}\{\omega \leq z\})^2 dz$	$\frac{1}{2m^2} \sum_{i=1}^m \sum_{k=1}^m  y_i - y_k  - \frac{1}{m} \sum_{i=1}^m  y_i - \omega $	b
Spherical	SS	$\frac{f_P(\omega)}{\left( \int_{-\infty}^{\infty} f_P^2(y) dy \right)^{1/2}}$	$f_P(\omega) \left[ \sum_{k=1}^m \left( \frac{f_P(y_k)}{\sum_{j=1}^m f_P(y_j)} \right)^2 \right]^{-1/2}$	a
Energy	ES	$\frac{1}{2} \mathbb{E}_P [ y - y^* ^\eta] - \mathbb{E}_P [ y - \omega ^\eta]$	$\frac{1}{2m^2} \sum_{i=1}^m \sum_{j=1}^m  y_i - y_j ^\eta - \frac{1}{m} \sum_{i=1}^m  y_i - \omega ^\eta$	b,c

<sup>a</sup>  $f_P(x)$  signifies the empirical density of  $P$  at  $x$ . Determined from eCDF using kernel smoothing

<sup>b</sup>  $y_i$  and  $y_j$  are independent draws from the distribution forecast  $P$

<sup>c</sup> Index  $\eta \in (0, 2)$ ; For  $\eta = 1$ , we yield  $S_{\text{CRPS}}(P, \omega)$  and  $\eta \rightarrow 2$  leads to  $S_{\text{SE}}(P, \omega) = -|\mu_P - \omega|^2$

### 6.1.4 Energy Score

940 Gneiting and Raftery (2007) proposed a generalization of the CRPS the so-called energy score

941

$$S_{\text{ES}}(P, \omega) = \frac{1}{2} \mathbb{E}_P [|y - y^*|^\eta] - \mathbb{E}_P [|y - \omega|^\eta], \quad (68)$$

942 where the index  $\eta \in (0, 2)$  and  $y$  and  $y^*$  are independent copies of the forecast distribution,  $P \in \mathcal{P}_\eta$ . This  
943 is a strictly proper score (Székely, 2003) and reduces to the CRPS for  $\eta = 1$  and the negative squared  
944 error  $S_{\text{SE}}(P, \omega) = -|\mu_P - \omega|^2$  in the limit of  $\eta \rightarrow 2$  (Gneiting and Raftery, 2007), where  $\mu_P$  is the mean  
945 of the distribution forecast  $P$ . For an ensemble forecast of  $m$  values  $\mathbf{y} = (y_1, \dots, y_m)^\top$ , the energy score  
946 may be computed as follows

947

$$S_{\text{ES}}(P, \omega) = \frac{1}{2m^2} \sum_{i=1}^m \sum_{j=1}^m |y_i - y_j|^\eta - \frac{1}{m} \sum_{i=1}^m |y_i - \omega|^\eta. \quad (69)$$

948 Table 6 summarizes the different *strictly proper* scoring rules for a distribution forecast  $P$  with density  
949 forecast  $p$  defined up to the  $\mu$ -measure zero. All of them reward leptokurtic forecast distributions with  
950 probability mass that centers on the verifying observation,  $\omega$ . As a reminder, the logarithmic score LS is  
951 the only tabulated scoring rule which ignores model predicted probabilities of all non-realized outcomes,  
952 thus, is strictly local only and highly sensitive to low probability events.

953 **6.2 Numerical Examples**

954 We use three different case studies to illustrate the computation of the scoring rules for univariate density  
955 forecasts. The first study provides a simple graphic illustration of the computation of the scoring rules.  
956 The second and third case study illustrate the application of scoring rules to discharge forecasts obtained  
957 from Bayesian Model Averaging and Generalized Likelihood Uncertainty Estimation.

958 **6.2.1 Case Study III: Graphical Illustration**

959 Suppose the forecast distribution  $P$  of the discharge  $\Omega$  in mm/d is exactly described by a gamma  
960 distribution  $P = \mathcal{G}(a_P, b_P)$  (see Figure 10) with dimensionless shape parameter  $a_P = 3$  and scale  
961 parameter  $b_P = 1$  mm/d. We slide the verifying observation  $\omega$  from left to right across the distribution  
962  $P$  and display the corresponding values of the quadratic (blue), logarithmic (red), spherical (green) and  
963 continuous ranked probability (yellow) scores of Table 6 on the  $y$ -axis. The different scoring rules exhibit  
964 a characteristic concave shape and provide the largest reward (smallest loss) in the high probability  
965 density region of the forecast distribution. Outside this region the different scoring rules decline in  
966 value with increasing distance from their maximum reward. The maximum of the quadratic, logarithmic  
967 and spherical scoring rules coincides exactly with the peak of the gamma distribution at  $\omega = 2$ . The  
968 maximum reward of the CRPS is well removed from the peak (mode) of the forecast distribution and  
969 concentrates on the median of  $P$  at about  $\omega = 2.67$ . The overall functional shape of the four scoring  
970 rules is rather similar but with appreciable differences in the curvature of QS, LS, SS and CRPS. The  
971 logarithmic score responds most strongly to changes in the density of the distribution forecast  $P$ , whereas  
972 the quadratic score displays a much more damped response to changes in the predictive density. Now  
973 imagine that  $\omega = 4.0$  mm/d is the discharge which materializes (vertical dashed line) at a future time.  
974 According to the gamma distribution forecast  $P = \mathcal{G}(3, 1)$  we yield  $f_{\mathcal{G}}(\omega, 3, 1) = 0.1465$  (black diamond)  
975 with values of the scoring rules (colored dots) equal to  $S_{QS} = 0.2912$ ,  $S_{LS} = -1.9206$  nats,  $S_{SS} = 3.375$   
976 and  $S_{CRPS} = -0.7585$  mm/d. As byproduct of another derivation, Appendix I presents a closed-form  
977 expression for the CRPS of a gamma distribution forecast. Equation (I.25) matches the yellow line.  
978 The strong similarities between the QS, SS and LS do not come as a surprise. They belong to a limited  
979 class of *strictly proper* scoring rules on  $\mathcal{L}_1(\Omega)$  and/or  $\mathcal{L}_2(\Omega)$  and, thus, are expected to be related. Indeed,  
980 under some regularity conditions, *Bernardo* (1979) has shown that every proper local scoring rule is

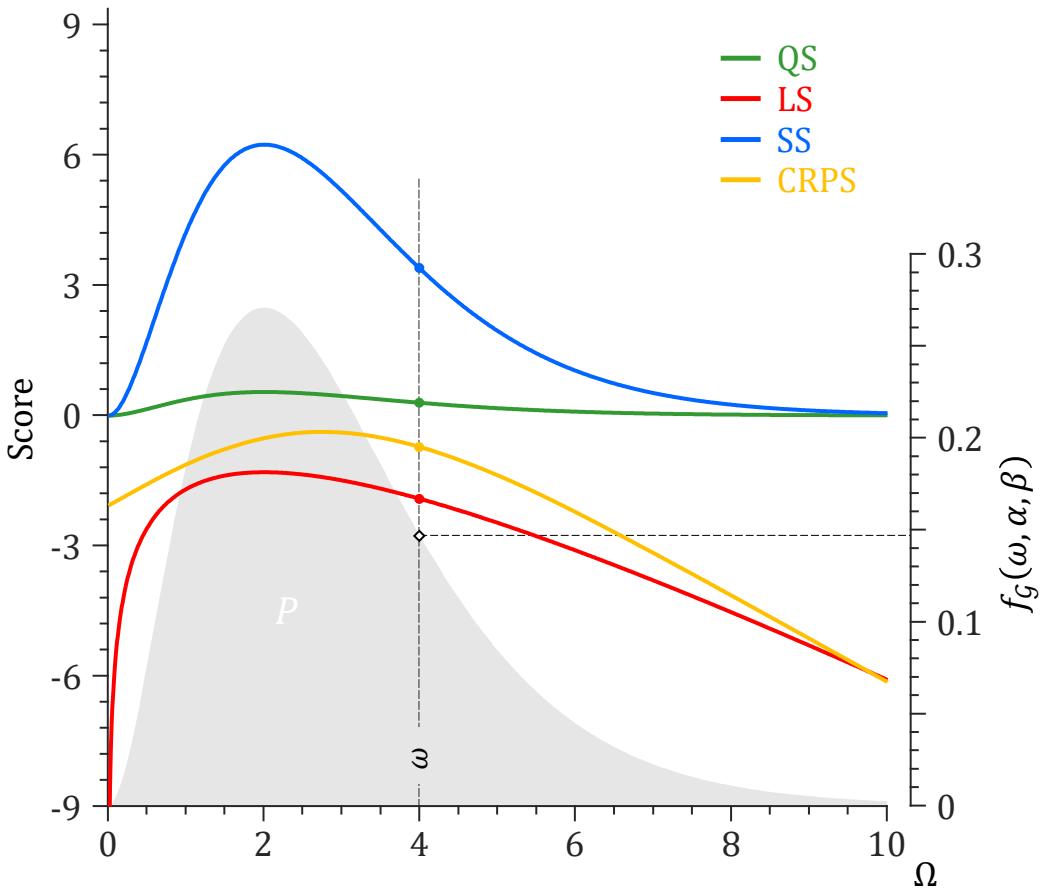


Figure 10: Hypothetical distribution forecast  $P$  of the discharge  $\Omega$  (mm/d) and traces of the quadratic (green), logarithmic (red), spherical (blue) and continuous ranked probability (yellow) scoring rules across the streamflow distribution computed using their numerical definitions in Table 6. The probabilistic forecast is a gamma distribution  $P = \mathcal{G}(a, b)$  with shape and scale parameters  $a = 3$  (-) and  $b = 1 \text{ mm/d}$ , respectively, and PDF  $f_G(\omega, a, b) = \Gamma^{-1}(a)b^{-a}\omega^{a-1} \exp(-\omega/b)$ , where  $\Gamma(x)$  is the gamma function.

981 equal to an affine transformation (45) of the logarithmic score. In principle, we only require one *strictly*  
982 *proper* scoring rule, though there are benefits in using multiple different scoring rules for verification of  
983 probabilistic model forecasts and/or multi-model ensemble prediction systems.

#### 984 6.2.2 Case Study IV: Simple Illustration

985 We revisit the distribution forecasts of Fig. 4a and compute the *strictly proper* scoring rules of Table  
986 6 using a collection of  $n = 10^4$  observations  $\omega_1, \dots, \omega_n$  drawn at random from the *true distribution*  $Q$ .  
987 Table 7 documents the outcome of this analysis and lists mean values of the quadratic, logarithmic,  
988 spherical and continuous ranked probability scores for  $P_1, \dots, P_5$ . Albeit different, these five distribution  
989 forecasts were deemed equivalent according to the coefficient of variation or  $C_v$  performance metric. The

Table 7: Mean values of the quadratic, logarithmic, spherical and continuous ranked probability scoring rules for the distribution forecasts  $P_1, \dots, P_5$  portrayed in Fig. 4a. The last column reports the mean scores for a perfect distribution forecast,  $P = Q$ .

Score	Fig. 4a					
	$P_1$ red	$P_2$ blue	$P_3$ green	$P_4$ yellow	$P_5$ purple	$P = Q$ gray
QS	0.484	0.467	0.501	0.513	0.525	0.534
LS	-6.520	-2.422	-2.190	-2.260	-2.479	-2.132
SS	6.197	6.148	6.390	6.332	6.454	6.608
CRPS	-0.649	-0.653	-0.663	-0.666	-0.642	-0.638

four scoring rules display a considerable variation among the distribution forecasts and assign different rewards to  $P_1, \dots, P_5$ . This testifies to their complete evaluation of the forecast distribution. The QS, LS, SS and CRPS assign their highest scores to  $P_3$  or the green forecast distribution in Fig. 4a. All four scoring rules agree that this forecast distribution is the best approximation of the *true distribution*  $Q$ , a conclusion that is supported by a visual comparison of  $P_1, \dots, P_5$  and  $Q$ . The last column lists the values of the scoring rules when the forecaster quotes  $P = Q$ . This is the maximum value each scoring rule can attain and equal to a Monte Carlo estimate of  $\mathcal{S}(Q, Q)$  in Equation (21). Thus, if we subtract each column  $P_1, \dots, P_5$  from this last column, we yield estimates of the scores divergences  $d_{\text{QS}}(P, Q)$ ,  $d_{\text{LS}}(P, Q)$ ,  $d_{\text{SS}}(P, Q)$  and  $d_{\text{CRPS}}(P, Q)$  for each distribution forecast.

### 6.2.3 Case Study V: Bayesian Model Averaging

We now illustrate the scoring rules by application to density forecasts of river discharge from a multi-model ensemble of  $K = 8$  conceptual hydrologic models of the Leaf River watershed ( $1950 \text{ km}^2$ ) located north of Collins, Mississippi. This ensemble is described in *Vrugt and Robinson (2007)* and interested readers are referred to this publication for more details. Figure 11 displays the discharge forecasts for a short but representative period of the training data set. The interval of the daily discharge forecasts generally envelops the measured discharge. Note that some models issue a negative forecast (days 1-40, 240-260 and 350-400) as a result of linear bias-correction as recommended by *Raftery et al. (2005)*.

Let  $\beta_k$ ,  $y_{kt}$  and  $f_k(y|y_{kt}, \psi_k)$  denote the weight, streamflow prediction and conditional density of the  $k$ th model of the ensemble at time  $t$ . The density of the multi-model forecast distribution  $P_t$  at  $t = (1, \dots, n)$

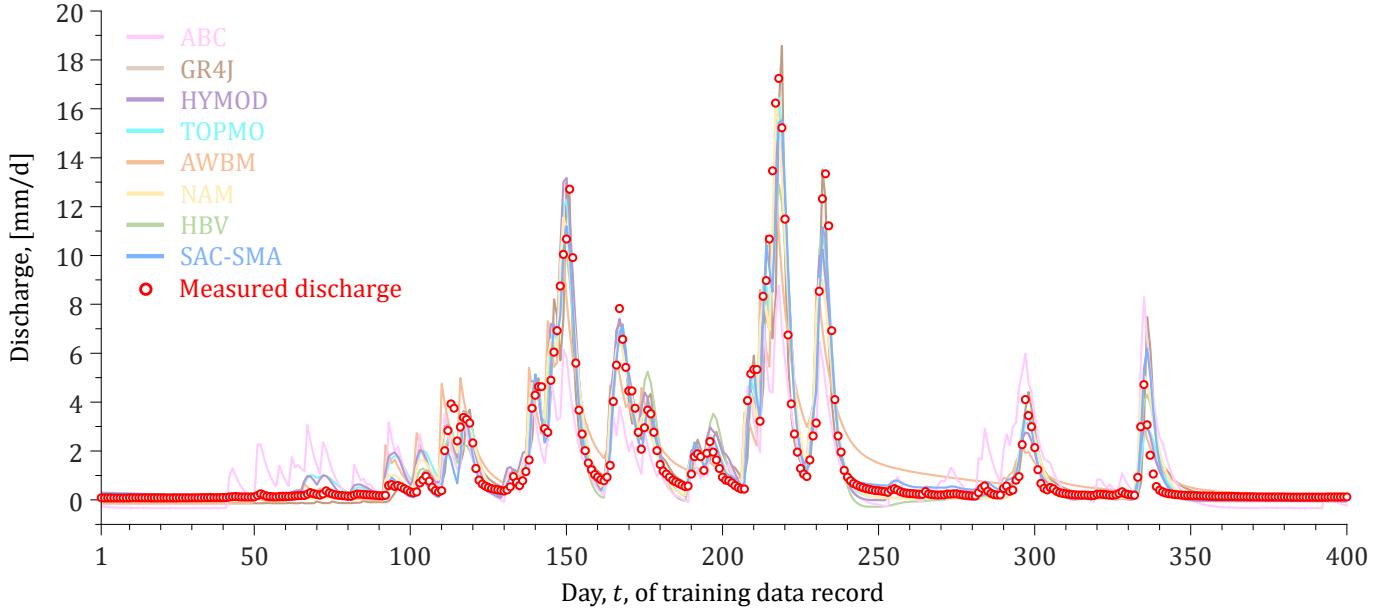


Figure 11: Streamflow predictions of the ABC (3), GR4J (4), HYMOD (5), TOPMO (6), AWBM (8), NAM (9), HBV (9) and SAC-SMA (13) conceptual watershed models for a small period of the  $n = 3,000$  day training data record. The number between parenthesis lists the number of calibration parameters. The solid red circles correspond to the daily measured discharges.

1009 is now equal to a mixture distribution of the models' conditional PDFs

$$\text{1010 } f_{P_t}(y|\boldsymbol{\beta}, \boldsymbol{\psi}) = \sum_{k=1}^K \beta_k f_k(y|y_{kt}, \psi_k), \quad (70)$$

1011 centered on the weighted-average forecast

$$\text{1012 } \mu_{P_t} = \sum_{k=1}^K \beta_k y_{kt} \quad (71)$$

1013 with weights,  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_K)^\top$ , constrained to the probability simplex,  $\Delta^{K-1} = \{\boldsymbol{\beta} \in \mathbb{R}^K : \beta_1 + \text{1014 } \dots + \beta_K = 1; \beta_k \geq 0 \text{ for } k = 1, \dots, K\}$  and shape parameters  $\psi_k$  of each model's predictive density 1015 assembled in the array  $\boldsymbol{\psi} = (\psi_1, \dots, \psi_K)^\top$ . Equations (70) and (71) are also known in the literature as 1016 the BMA forecast density and BMA model forecast, respectively (*Raftery et al.*, 2005). Table 8 presents 1017 mathematical formulations of the conditional PDFs,  $f_k(y|y_{kt}, \psi_k)$ ;  $k = (1, \dots, K)$ , used to construct the 1018 forecast density.

Table 8: Formulation, coefficients and unknown parameters of the 1. generalized normal, 2. lognormal, 3. gamma and 4. Weibull predictive PDFs used in the BMA forecast density of Equation (70).

Formulation	Coefficients	Unknown parameters
1. $f_k(y \mu_{kt}, s_k^2, \tau_k) = \frac{\tau_k}{2s_k\Gamma(\tau_k^{-1})} \exp\left[-\left(\frac{ y - \mu_{kt} }{s_k}\right)^{\tau_k}\right]$	$\mu_{kt} = y_{kt}$	$\Psi = (s_1^2, \dots, s_K^2, \tau_1, \dots, \tau_K)^\top \S\dagger$
2. $f_k(y \mu_{kt}, s_k^2) = \frac{1}{ys_k\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{\log(y) - \mu_{kt}}{s_k}\right)^2\right]$	$\mu_{kt} = \log(y_{kt}) - \frac{1}{2}s_k^2$	$\Psi = (s_1^2, \dots, s_K^2)^\top \S*$
3. $f_k(y a_{kt}, b_{kt}) = \frac{1}{\Gamma(a_{kt})b_{kt}^{a_{kt}}} y^{a_{kt}-1} \exp\left(-\frac{y}{b_{kt}}\right)$	$a_{kt} = \frac{ y_{kt} ^2}{s_k^2}, b_{kt} = \frac{s_k^2}{ y_{kt} }$	$\Psi = (s_1^2, \dots, s_K^2)^\top \S\dagger$
4. $f_k(y \lambda_{kt}, s_k^2) = \frac{s_k}{\lambda_{kt}} \left(\frac{y}{\lambda_{kt}}\right)^{s_k-1} \exp\left[-\left(\frac{y}{\lambda_{kt}}\right)^{s_k}\right]$	$\lambda_{kt} = \frac{ y_{kt} }{\Gamma(1+1/s_k)}$	$\Psi = (s_1^2, \dots, s_K^2)^\top \ \P$

$\S$  We use a (i) group,  $s_1^2 = \dots = s_K^2$ , and (ii) sole,  $s_1^2, \dots, s_K^2$ , variance for the  $K = 8$  models

$\dagger$  The variance  $\sigma_k^2$  of  $f_k(y|\cdot)$  is equal to  $s_k^2$

$*$  The variance of the lognormal density is equal to  $\sigma_k^2 = (\exp(s_k^2) - 1) \exp(2\mu_{kt} + s_k^2)$

$\|$  We use  $s_k^2$  for the shape parameter, thus, evaluate a group and model dependent (= sole) shape parameter

$\P$  The variance of the Weibull density is equal to  $\sigma_k^2 = \lambda_{kt}^2 [\Gamma(1+2/s_k) - \Gamma^2(1+1/s_k)]$

1019 The coefficients of the conditional PDFs are defined so that the means of the generalized normal, lognormal,  
1020 gamma and Weibull distributions coincide with the respective model forecasts,  $y_{kt}$ , where  $k = (1, \dots, K)$ .

1021 The shape parameter  $\tau_k$  determines the kurtosis of the generalized normal density. A value of  $\tau_k = 2$   
1022 results in a normal distribution (albeit with variance  $\sigma_k^2/2$ ), a value of  $\tau_k = 1$  equates to a Laplace  
1023 (double-exponential) distribution, and  $\tau_k \rightarrow \infty$  converges to a uniform density on  $[y_{kt} - \sigma_k, y_{kt} + \sigma_k]$  and  
1024 a zero density outside this interval. Thus, the larger the value of  $\tau_k$  the less peaked the conditional PDF  
1025 of the  $k$ th model will be and the more tight its associated prediction intervals around  $y_{kt}$ .

1026 If we assume that the models' forecast errors are independent, then the  $d$ -vector  $\Theta = (\beta, \Psi)$  of weights  $\beta$   
1027 and shape parameters  $\Psi$  of the conditional PDFs of Table 8 can be determined from maximization of the  
1028 BMA log-likelihood function,  $\ell(\beta, \Psi | \omega) = \sum_{t=1}^n \log(f_{P_t}(\omega_t | \beta, \Psi))$ , using MCMC simulation with the  
1029 DREAM algorithm (*Vrugt et al.*, 2008) and weights constrained to the probability simplex. Although the  
1030 model ensemble does not satisfy the independence assumption, this should not affect much the estimates  
1031 of the weights  $\beta$  and shape parameters  $\Psi$ , because we are estimating the conditional distribution for a  
1032 scalar observation given  $K$  forecasts, rather than for several observations simultaneously (*Raftery et al.*,  
1033 2005).

1034 The variance of the BMA forecast density  $f_{P_t}(y|\beta, \Psi)$  in Equation (70) is equal to

$$\sigma_{P_t}^2 = \sum_{k=1}^K \beta_k (\sigma_k^2 + y_{kt}^2) - \mu_{P_t}^2 \quad (72)$$

1036 where  $y_{kt}$  and  $\sigma_k^2$  denote the mean and variance of the  $f_k(y|y_{kt}, \psi_k)$ 's at time  $t$ . As the conditional PDFs  
 1037 of Table 8 admit a closed-form solution for their variances,  $\sigma_k^2$ ;  $k = (1, \dots, K)$ , the coefficient of variation  
 1038 of the BMA forecast density at time  $t$  is exactly defined;  $C_{v,t} = \sigma_{P_t}/\mu_{P_t}$ . Lower and upper endpoints of  
 1039 the  $\gamma = 100(1 - \alpha)\%$  prediction interval of the BMA mixture density can be derived from the CDF

$$F_{P_t}(y|\boldsymbol{\beta}, \boldsymbol{\Psi}) = \sum_{k=1}^K \beta_k F_k(y|y_{kt}, \psi_k), \quad (73)$$

1040 so that  $F_{P_t}(l_t|\boldsymbol{\beta}, \boldsymbol{\Psi}) = \alpha/2$  and  $F_{P_t}(u_t|\boldsymbol{\beta}, \boldsymbol{\Psi}) = 1 - \alpha/2$ . At each time  $t$ , we solve for the lower and upper  
 1041 predictive quantiles at different  $\alpha$  values using an iterative root finding procedure. If we evaluate the  
 1042 CDF in Equation (73) at each verifying observation and sort the resulting values in ascending order, then  
 1043 the reliability  $R_l$  of the BMA forecast distribution is easily computed using Equation (19). We use the  
 1044 MODELAVG toolbox of *Vrugt* (2018) in MATLAB to determine maximum likelihood values of the BMA  
 1045 model parameters for the conditional PDFs of Table 8 along with performance metrics, scoring rules and  
 1046 discharge prediction intervals of the BMA forecast density.

1047 Figure 12 presents traces of the weighted-average BMA forecast (black line) and associated 50, 75 and 95%  
 1048 prediction intervals (gray regions) using the (a) normal ( $\tau = 2$ ), (b) lognormal, (c) generalized normal, (d)  
 1049 gamma and (e) Weibull predictive PDFs of Table 8 with a sole forecast variance,  $\sigma_k^2$ ;  $k = (1, \dots, K)$  for  
 1050 each watershed model. The solid red squares depict the measured daily discharges. This weighted-average  
 1051 BMA forecasts describe the discharge observations quite well and appear relatively unaffected by the  
 1052 choice of the models' predictive distribution. This does not hold, however for the BMA forecast density.  
 1053 Indeed, we observe large differences in the spread of the 50, 75 and 95% prediction intervals of the BMA  
 1054 mixture distribution of Equation (70) among the conditional PDFs of Table 8. These differences are  
 1055 particularly well visible at the higher flows, for example between days 2,740 - 2,750, but also manifest  
 1056 themselves at lower flows, particularly during the dry period between days 2,800 and 2,870 of the training  
 1057 record. The symmetry of the normal and generalized normal distributions translates into BMA prediction  
 1058 interval with an approximately constant spread, independent of flow level. As a result, the prediction  
 1059 uncertainty of the BMA model is relatively small for the driven part of the hydrograph and comparatively  
 1060 large at the lower flows. The lognormal, gamma and Weibull distributions are defined only for  $y \geq 0$  and,  
 1061 thus, are expected to yield a more accurate and fluent description of streamflow behavior. Indeed, the  
 1062 BMA prediction intervals of these three skewed distributions are comparatively large at the highest flows  
 1063 but their spread declines rapidly with discharge magnitude and collapses to a very small region around

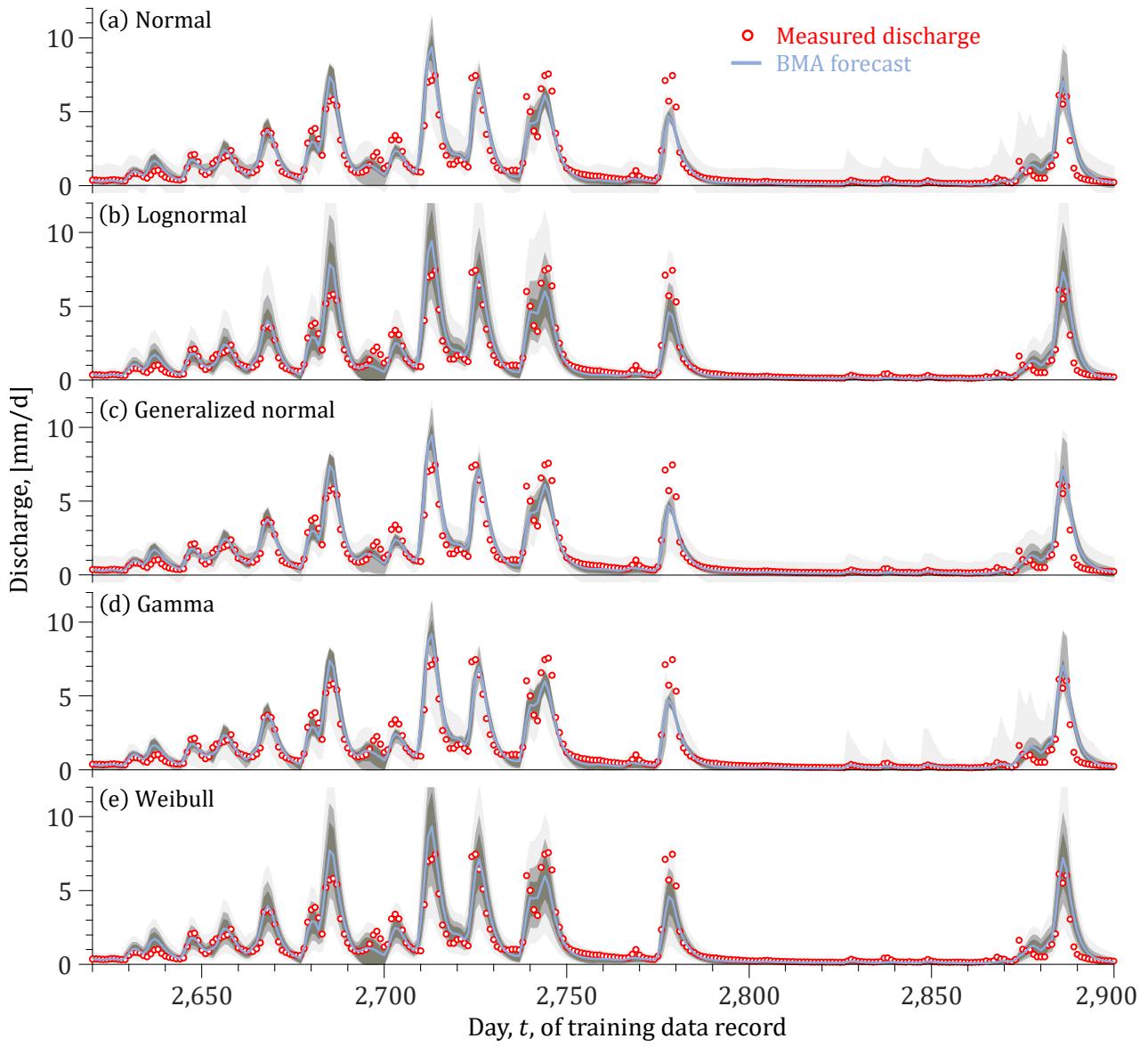


Figure 12: Weighted average BMA forecast (solid blue line) and 50% (dim gray), 75% (medium gray) and 95% (light gray) quantiles of the BMA forecast distribution for a representative 270-day period of the training data set using the (a) normal, (b) lognormal, (c) generalized normal, (d) gamma and (e) Weibull distributions with a model-dependent (= sole) forecast variance. The red circles are the daily discharge observations.

the mean forecast (solid blue line) during baseflow. The gamma distribution produces much sharper BMA density forests during peak flows than the lognormal and Weibull distributions but this is at the expense of an overly dispersed BMA forecast distribution for small rainfall events in the long dry period between days 2,800 and 2,870. This will deteriorate the average spread and performance of the gamma distribution. The normal and generalized normal distributions display similar jagged BMA prediction intervals in the long recession period.

The BMA forecast density offers an excellent opportunity for application of the scoring rules. For each

1072 scoring rule, we compute a (time-averaged) mean score

$$1073 \quad \bar{S}_{\text{XX}}(\mathbf{P}, \boldsymbol{\omega}) = \frac{1}{n} \sum_{t=1}^n S_{\text{XX}}(P_t, \omega_t), \quad (74)$$

1074 where  $\mathbf{P} = \{P_1, \dots, P_n\}$  is the collection of probabilistic forecasts derived from the BMA model. Table  
 1075 9 reports the mean values of the quadratic, logarithmic, spherical, continuous ranked probability and  
 1076 energy scoring rules for the BMA forecast density of Equation (70) with a normal, lognormal, generalized  
 1077 normal, gamma and Weibull predictive PDF using a constant *group* or *sole* forecast variance. We also  
 1078 list the values of the performance metrics of Table 1 and report the log-likelihood

$$1079 \quad \ell(\boldsymbol{\beta}, \boldsymbol{\psi} | \boldsymbol{\omega}) = \sum_{t=1}^n \left\{ \log_e \left( \sum_{k=1}^K \beta_k f_k(\omega_t | y_{kt}, \psi_k) \right) \right\}, \quad (75)$$

1080 Root Mean Square Error (RMSE)

$$1081 \quad s_{\text{RMSE}}(\boldsymbol{\mu}_P, \boldsymbol{\omega}) = \sqrt{\frac{1}{n} \sum_{t=1}^n (\omega_t - \mu_{P_t})^2}, \quad (76)$$

1082 coefficient of determination ( $R^2$ )

$$1083 \quad s_{R^2}(\boldsymbol{\mu}_P, \boldsymbol{\omega}) = 1 - \frac{\sum_{t=1}^n (\omega_t - \mu_{P_t})^2}{\sum_{t=1}^n (\omega_t - m_\omega)^2}, \quad (77)$$

1084 and KG efficiency

$$1085 \quad s_{\text{KGE}}(\boldsymbol{\mu}_P, \boldsymbol{\omega}) = 1 - \sqrt{(a-1)^2 + (b-1)^2 + (r-1)^2}, \quad (78)$$

1086 of the weighted-average BMA forecast of Equation (71) where  $a = m_{\mu_P}/m_\omega$  and  $b = s_{\mu_P}/s_\omega$  are the unitless  
 1087 ratios of the sample means and sample standard deviations, respectively, of the weighted-average BMA  
 1088 discharge forecasts,  $\boldsymbol{\mu}_P = (\mu_{P_1}, \dots, \mu_{P_n})^\top$ , and verifying measurements,  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_n)^\top$ , respectively,  
 1089 and  $r$  is the sample Pearson correlation coefficient of the  $n$ -data pairs,  $\{(\omega_1, \mu_{P_1}), \dots, (\omega_n, \mu_{P_n})\}$ . The  
 1090 tabulated data highlight several important findings.

1091 1. The time-averaged values of the *strictly proper* scoring rules of the BMA forecast density differ  
 1092 substantially between the predictive PDFs of Table 8. This is particularly true for the QS, LS and  
 1093 SS, whereas the CRPS and ES demonstrate considerably less variation.

1094 2. The use of a model-specific forecast variance generally increases the values of the scoring rules for  
 1095 the BMA forecast density but this improvement is marginal for the lognormal, generalized normal  
 1096 and Weibull distributions.

Table 9: Time-averaged values of the *strictly proper* scoring rules of Table 6 for the BMA density forecast  $f_{P_t}(y|\boldsymbol{\beta}, \boldsymbol{\Psi})$ ;  $t = (1, \dots, n)$  of Equation (70) using the normal ( $\tau = 2$ ), lognormal, generalized normal, gamma and Weibull predictive PDFs with a group or model-dependent (= sole) constant variance. We also list the performance metrics,  $R_l$ ,  $C_v$ ,  $C$  and  $W$  of Table 1 and report the log-likelihood,  $\ell(\boldsymbol{\beta}, \boldsymbol{\Psi}|\boldsymbol{\omega})$ , RMSE,  $R^2$  and KG efficiency of the weighted-average BMA forecast. The bottom row lists the number  $d$  of unknown BMA parameters.

	Normal		Lognormal		Gen. normal		Gamma		Weibull	
	group	sole	group	sole	group	sole	group	sole	group	sole
QS	1.345	2.862	3.697	3.732	2.864	2.922	0.995	2.657	3.603	3.616
LS	-1.162	-0.291	0.145	0.153	-0.265	-0.257	-1.426	-0.408	0.103	0.108
SS	15.28	7.695	33.85	29.79	10.33	11.87	8.018	11.13	36.91	36.97
CRPS	-0.265	-0.238	-0.220	-0.220	-0.236	-0.236	-0.261	-0.237	-0.221	-0.220
ES	-0.502	-0.519	-0.513	-0.525	-0.516	-0.514	-0.563	-0.528	-0.520	-0.520
$R_l$	0.229	0.105	0.135	0.128	0.104	0.098	0.234	0.128	0.129	0.127
$C_v$	0.789	-3.917	0.772	0.859	1.849	1.784	2.189	-2.111	0.800	0.867
$C$	0.961	0.962	0.972	0.971	0.951	0.951	0.958	0.956	0.967	0.969
$W$	2.214	2.557	1.843	1.829	2.015	2.047	1.959	1.513	1.661	1.698
$\ell(\boldsymbol{\theta} \boldsymbol{\omega})$	-2416	-604.8	300.7	318.0	-550.1	-534.5	-2966	-847.8	215.2	223.8
RMSE	0.708	0.716	0.721	0.724	0.719	0.715	0.750	0.726	0.722	0.725
$R^2$	0.866	0.863	0.861	0.860	0.862	0.864	0.850	0.860	0.861	0.860
KGE	0.868	0.863	0.869	0.865	0.862	0.865	0.824	0.854	0.867	0.862
$d$	9	16	9	16	17	24	9	16	9	16

- 1097 3. The lognormal and Weibull distributions maximize the overall performance of the BMA forecast  
 1098 density according to the *strictly formal* scoring rules. This confirms the advantages of using a  
 1099 conditional PDF commensurate with the skewed distribution of the discharge data. Perhaps rather  
 1100 surprisingly, but this advantage does not hold for the gamma distribution. While the spread of  
 1101 the BMA forecast density of this distribution is much smaller on average than the lognormal and  
 1102 Weibull distributions, the scoring rules of the gamma distribution suffer primarily on days with an  
 1103 underdispersed BMA forecast density (e.g. see days 2,740-2,750 and around 2,780) and to a lesser  
 1104 extent on days with overly dispersed BMA prediction intervals in the dry period (see Fig. 12d).
- 1105 4. The reliability,  $R_l$ , varies quite a bit among the different conditional PDFs of the BMA forecast  
 1106 density. This measure of the statistical consistency of the BMA distribution forecasts and the  
 1107 observed discharge data does not prefer the lognormal and Weibull distributions but rather assigns  
 1108 the lowest value of about 0.10 to the generalized normal distribution. But as was illustrated in Fig.  
 1109 4, the reliability is not a sufficient condition, hence, may support an erroneous ranking of the BMA

1110 models.

- 1111 5. The coefficient of variation  $C_v$  exhibits a large variation among the different conditional PDFs of  
1112 the BMA forecast density with values ranging between -3.9 and 2.2 for the normal and gamma  
1113 distributions, respectively. Time-averaging obfuscates interpretation of this measure of relative  
1114 dispersion of the BMA forecast distribution as its value may be corrupted by large  $C_v$  values on  
1115 days with near-zero discharges. This may at least explain in part why the  $C_v$  correlates poorly with  
1116 the average spread  $W$  of the 95% BMA prediction intervals. The listed values of the  $C_v$  are exact  
1117 for the BMA mixture density of Equation (70), but are otherwise difficult to estimate accurately,  
1118 particularly for skewed discharge distributions with a mean close to zero.
- 1119 6. All distributions of Table 8 achieve an adequate coverage of  $C \approx 0.95$  at significance level  $\alpha = 0.05$ .  
1120 This testifies to the appropriate form of the likelihood function and the DREAM algorithm's ability  
1121 in successfully locating the maximum likelihood values of the BMA weights and shape parameters.  
1122 The coverage, however, provides an incomplete evaluation of density forecasts as pointed out in  
1123 Figure 4. We can turn the coverage into a *proper* scoring rule as will be discussed in the context of  
1124 quantile prediction.
- 1125 7. The tabulated data highlight the limitations of the conjectured sharpness principle of *Gneiting et al.* (2007). The gamma distribution with a model-dependent (sole) forecast variance results in  
1126 the smallest spread (width) of the 95% BMA prediction intervals, yet, is inferior to almost all other  
1127 conditional PDFs according to the *strictly proper* scoring rules and log-likelihood  $\ell(\boldsymbol{\beta}, \boldsymbol{\psi} | \mathbf{w})$  of the  
1128 time-averaged BMA forecast. Sharpness is a measure of the concentration (dispersion) of the BMA  
1129 forecast density, hence is an *improper* scoring rule as it is a property of the predictive distribution  
1130 only.
- 1132 8. The RMSE,  $R^2$  and KG efficiency of the weighted-average BMA forecast of Equation (71) appear  
1133 insensitive to the choice of forecast variance (group or sole) and conditional PDF. These summary  
1134 metrics are unresponsive to the large differences in the BMA forecast density shown in Fig. 12  
1135 between the predictive PDFs of Table 8. The log-likelihood on the contrary is much more responsive  
1136 to changes in the BMA forecast distribution with values of  $\ell(\boldsymbol{\beta}, \boldsymbol{\psi} | \mathbf{w})$  ranging between values of  
1137 -2,966 and 318 for the different conditional PDFs. The log-likelihood demonstrates the advantages

1138 of using a model-dependent (sole) forecast variance and confirms the superiority of the lognormal  
 1139 and Weibull distributions. In fact, the log-likelihood exhibits a perfect 1 : 1 relationship with the  
 1140 logarithmic scoring rule and an almost perfect linear relationship with the QS and CRPS (see Table  
 1141 10).

Table 10: Pearson correlation coefficients of the log-likelihood  $\ell(\boldsymbol{\beta}, \boldsymbol{\psi} | \boldsymbol{\omega})$  of the weighted-average BMA forecast and time-averaged values of the QS, LS, SS, CRPS and ES *strictly proper* scoring rules, improper performance metrics of Table 1 and RMSE,  $R^2$  and KGE scoring functions.

	QS	LS	SS	CRPS	ES	$R_l$	$C_v$	$C$	$W$	RMSE	$R^2$	KGE
$\ell(\boldsymbol{\theta}   \boldsymbol{\omega})$	0.998	1.000	0.690	0.976	0.361	0.809	0.596	0.560	-0.394	-0.351	0.360	0.666

1142 Thus, the blatant insensitivity of the RMSE,  $R^2$  and KG efficiency to the BMA forecast distribution  
 1143 is not a necessary consequence of these metrics acting only on the weighted-average BMA forecast.

1144 We have learned that widely used metrics such as the RMSE,  $R^2$  and KGE are woefully inadequate for  
 1145 evaluating distribution forecasts. This is a concerning finding as these metrics are commonly used by  
 1146 researchers and practitioners to evaluate model performance. The *strictly proper* scoring rules equip the  
 1147 hydrologist with an arsenal of "measures" for robustly evaluating distribution forecasts.

1148 To provide insights into the temporal dynamics of the scoring rules, please consider Figure 13 which  
 1149 presents a time series plot of the (a) BMA forecast distribution and traces of the quadratic (green),  
 1150 logarithmic (red), spherical (blue), continuous ranked probability (yellow) and energy (cyan) scoring  
 1151 rules. The scoring rules vary dynamically across the hydrograph and display a rapid succession of small  
 1152 and occasionally larger fluctuations in a pattern independent of flow level. We briefly summarize the  
 1153 most important differences between the traces of the scoring rules. The LS appears most responsive to  
 1154 the BMA forecast distribution with relatively large day-to-day variations in the value of this scoring  
 1155 rule. The trace of the QS is almost similar to that of the LS but with a noticeably smaller amplitude  
 1156 during the first 150 days of the displayed record. The strong agreement in the temporal behavior of these  
 1157 two scoring rules confirms the similarities in the mathematical definition of the LS and QS. The SS is  
 1158 seemingly unresponsive to the rapid succession of storm events in the first half of the record as a result  
 1159 of the large values this scoring rule attains towards the end of the so-called dry period. The CRPS and  
 1160 to a lesser extent the ES show most variation during the wet period when streamflow is intermediate to  
 1161 high but are almost unresponsive during periods without rainfall when discharge is low.

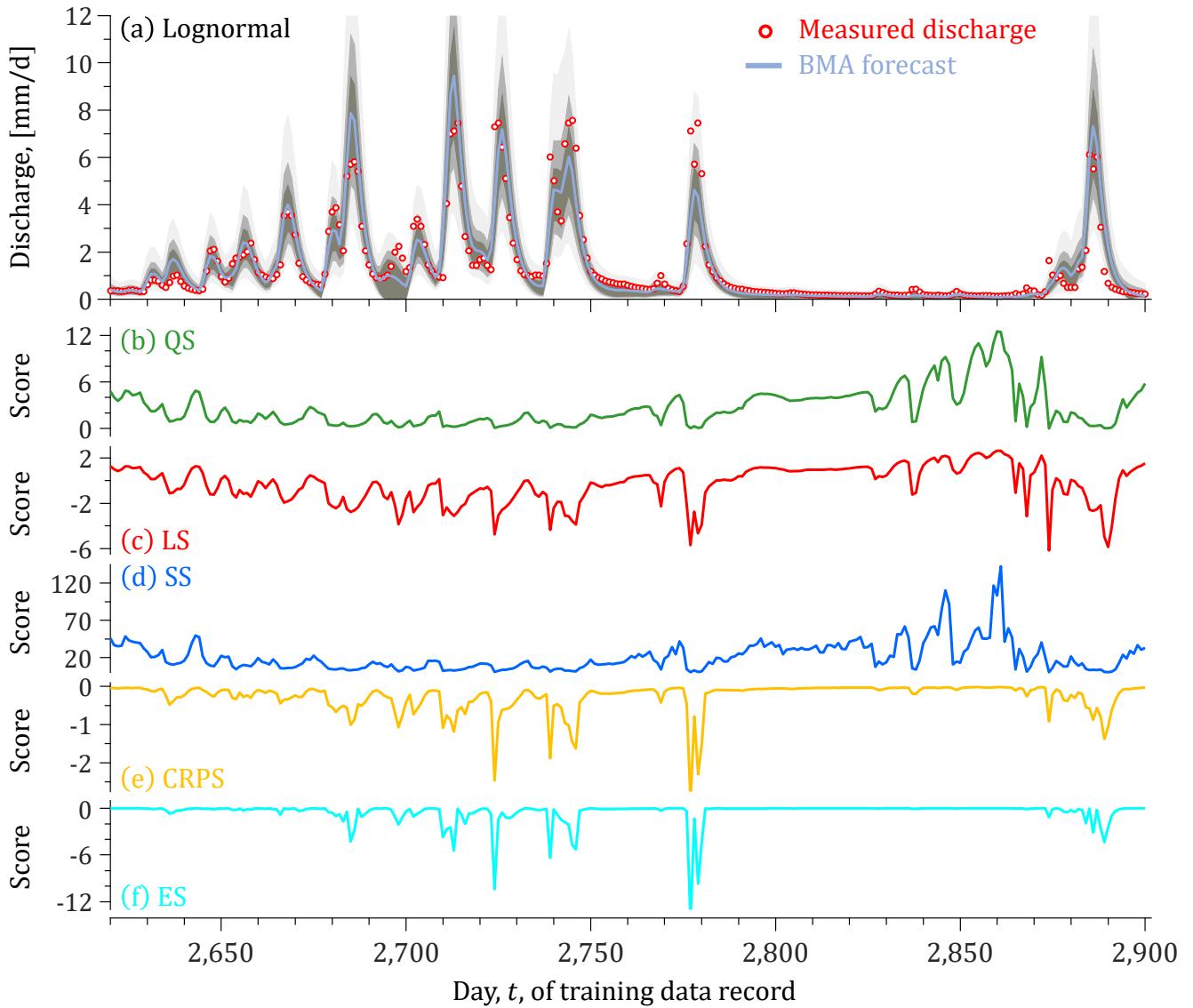


Figure 13: (a) 50%, 75% and 95% prediction intervals of the BMA forecast density for the lognormal distribution of Table 8 with a constant model-dependent (= sole) forecast variance and corresponding traces of the (b) quadratic, (c) logarithmic, (d) spherical, (e) continuous ranked probability and (f) energy scoring rules.

1162 The scoring rules show important similarities and differences. The QS, LS and SS attain their largest  
 1163 values at the lowest flow levels when the discharge observation is contained within the high probability  
 1164 density (HPD) region of the BMA forecast density and this distribution is relatively tight. This comes at  
 1165 no surprise as a narrow forecast distribution will return, on average, higher probability densities than a  
 1166 much more dispersed forecast distribution. The CRPS and ES respond quite differently and do not peak  
 1167 when the BMA forecast distribution is accurate and narrow but instead retain values near zero. What  
 1168 all scoring rules have in common is that they achieve the lowest values when the discharge measurements  
 1169 materialize outside the HPD region and in the tails of the BMA forecast distribution. In principle, any

of the *strictly proper* rules will suffice in evaluating the quality of the BMA forecast distribution and ranking the conditional PDFs. But the question which one(s) to use remains largely open (*Gneiting and Raftery, 2007; Alexander et al., 2022*).

#### 6.2.4 Case Study VI: Generalized Likelihood Uncertainty Estimation

Before we move on to scoring rules for multi-variable forecasts, we illustrate the application of scoring rules for model training and evaluation using application of the Generalized Likelihood Uncertainty Estimation (GLUE) method of *Beven and Binley* (1992) to the HYdrologic MODEL (HYMOD, *Boyle 2001*), Hydrologic model (Hmodel, *Schoups et al.* 2010) and Sacramento Soil Moisture Accounting (SAC-SMA) model (*Burnash et al.*, 1973). Details of these models can be found in the cited references and numerous other publications. Readers are referred to Appendix H for a schematic description of each model along with mathematical formulations of the fluxes between the control volumes and a table with unknown parameters. We simulate the rainfall-discharge relationship for the 3,000-day training record of the Leaf River watershed using daily estimates of areal average rainfall and potential evapotranspiration. The model equations are solved using a mass-conservative second-order integration method with adaptive time stepping. This guarantees a robust and accurate numerical solution of the simulated fluxes and state variables. A one-year spin-up period mitigates the impact of state variable initialization.

We draw at random  $m = 25,000$  parameter vectors  $\boldsymbol{\Theta}_1, \dots, \boldsymbol{\Theta}_m$  from the models' prior parameter ranges using Latin hypercube sampling. Next, we evaluate each parameter vector,  $\boldsymbol{\Theta}_i = (\theta_{i1}, \dots, \theta_{id})^\top$  and compute the quality-of-fit of its simulated daily discharge record,  $\mathbf{y}(\boldsymbol{\Theta}_i) = (y_1(\boldsymbol{\Theta}_i), \dots, y_n(\boldsymbol{\Theta}_i))^\top$ , using the normal log-likelihood

$$\ell(\boldsymbol{\Theta}_i | \boldsymbol{\omega}) = -\frac{n}{2} \log \left( \sum_{t=1}^n (\omega_t - y_t(\boldsymbol{\Theta}_i))^2 \right), \quad (79)$$

where  $e_t(\boldsymbol{\Theta}) = \omega_t - y_t(\boldsymbol{\Theta})$  signifies the  $t$ th residual. Then, we sort the simulations of the parameter vectors in decreasing likelihood. The result is a  $n \times m$  matrix  $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_m\}$  where the  $n \times 1$ -vector  $\mathbf{y}_1 = (y_{11}, \dots, y_{n1})^\top$  maximizes the log-likelihood of Equation (79) and  $\mathbf{y}_m = (y_{1m}, \dots, y_{nm})^\top$  returns the smallest value of  $\ell(\boldsymbol{\Theta} | \boldsymbol{\omega})$ . Then, for any sufficiently large ensemble size the simulated values  $y_{t1}, \dots, y_{tK}$  at time  $t$  make up a discrete forecast distribution  $P_t$  whose PDF  $f_{P_t}$  may be approximated using a normal kernel function with bandwidth proportional to the number of points  $K$ .

Figure 14 displays the relationship between the ensemble size  $K$  and the time-averaged value of the

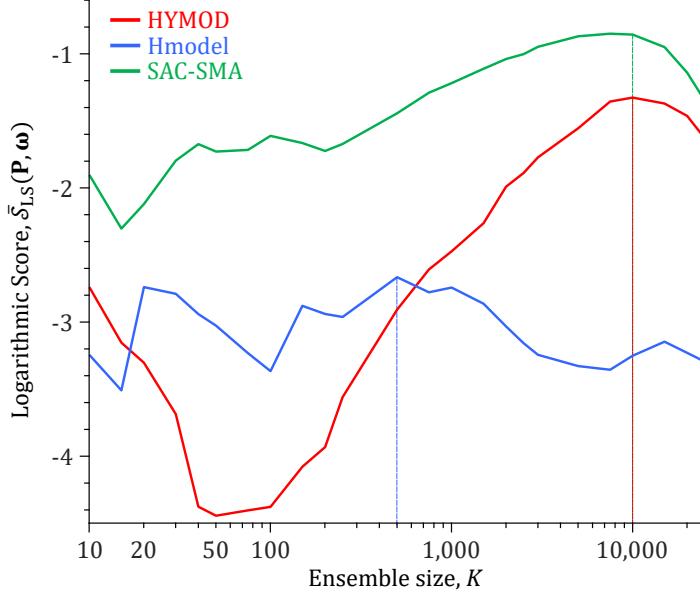


Figure 14: Relationship between the ensemble size,  $K$ , and the time-averaged value of the logarithmic scoring rule for HYMOD (red line), Hmodel (blue line) and SAC-SMA (green line) using daily discharge data from the Leaf River watershed. The dashed vertical lines point out the optimum size of the ensemble.

logarithmic scoring rule,  $\bar{S}_{LS}(\mathbf{P}, \boldsymbol{\omega})$ . As a reminder, an ensemble of size  $K$  uses the first  $K$  simulations  $\{\mathbf{y}_1, \dots, \mathbf{y}_K\}$  of  $\mathbf{Y}$ . The three models display quite a different  $\bar{S}_{LS}(K)$  relationship. The logarithmic score of HYMOD and the SAC-SMA model displays an obvious dependence on the ensemble size, while the  $\bar{S}_{LS}(K)$  curve of the Hmodel meanders up and down the graph in a much more unpredictable manner. The logarithmic score is comparatively small for small ensemble sizes but steadily increases to attain a maximum value at  $K = 500$  (Hmodel) and  $K = 10,000$  for HYMOD and the SAC-SMA model. Beyond this optimum the value of the LS degrades either steadily (HYMOD and SAC-SMA) or intermittently (Hmodel) with ensemble size. Secondly, the optimum ensemble size of Hmodel is substantially smaller than that of HYMOD and the SAC-SMA model. This could have several causes. The default parameter ranges of Hmodel may be too wide for the Leaf River resulting in a proportionally large number of inferior simulations in the GLUE sample whose inclusion in the ensemble can only deteriorate the forecast distribution. But the small ensemble of Hmodel may also signal model inadequacy. Thirdly, the SAC-SMA model consistently achieves the largest values of the LS followed by HYMOD and Hmodel. The optimum ensemble of HYMOD attains a substantially larger LS value than Hmodel, which in turn has two more adjustable parameters. The structure and/or process descriptions of the 7-parameter Hmodel (Fig. H.2) can only be inferior to that of the 5-parameter HYMOD (Fig. H.1) in describing the rainfall-discharge relationship of the Leaf River. Note that the  $\bar{S}_{LS}(K)$  curves hardly differ between

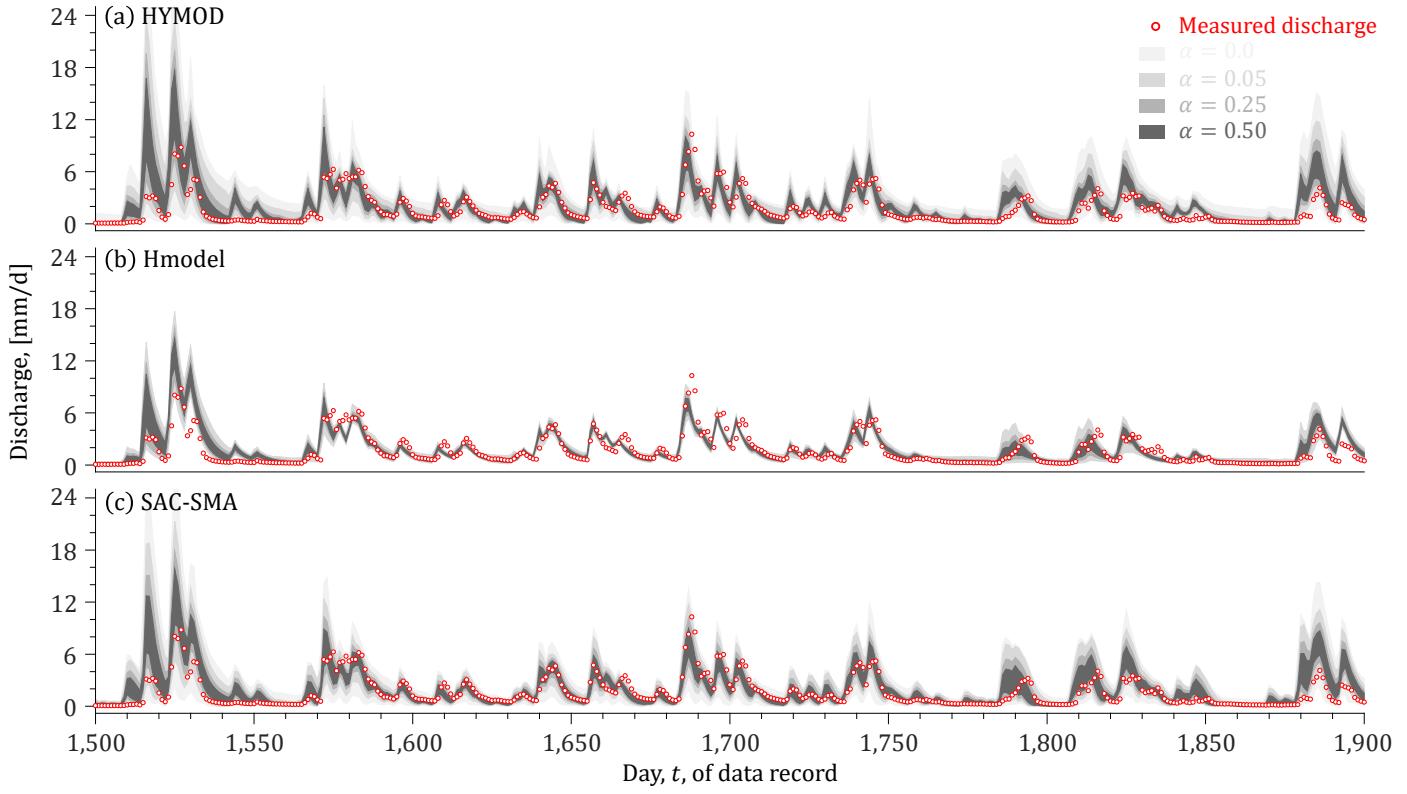


Figure 15: The 50, 75, 95 and 100% simulation uncertainty intervals of the (a) HYMOD, (b) Hmodel and (c) SAC-SMA conceptual hydrologic models for the optimum ensemble according to the logarithmic scoring rule.

1215 different GLUE trials.

1216 Next, we display the simulation uncertainty of the optimum ensemble (Figure 15) for (a) HYMOD, (b)  
1217 Hmodel and (c) the SAC-SMA model. Table 11 reports different summary statistics of the so-obtained  
1218 forecast distributions. The three graphs demonstrate several important findings.

1219 1. The different models describe the discharge observations quite well but with the exception of some  
1220 storm events, for example, on days 1,510 - 1,530 and 1,880 - 1,900 of the training data record.

1221 All three models significantly overestimate the measured discharge. Upon visual inspection of the  
1222 simulated discharge records, we conclude that this discrepancy is not a consequence of an inaccurate  
1223 initial state prior to the onset of the storm event but rather signals faulty rainfall estimates.

1224 2. The forecast distributions of HYMOD and SAC-SMA appear quite similar with intervals of the  
1225 simulated discharge that are in close agreement with each other and much wider on average than  
1226 those obtained from Hmodel. This is confirmed by the tabulated values of the width  $W$  of the  
1227 95% prediction intervals, which equal 2.401, 1.549 and 2.386 mm/d for HYMOD, Hmodel and

1228 SAC-SMA, respectively.

- 1229 3. The relatively sharp forecast distributions of Hmodel are at the expense of an insufficient coverage  
 1230 of its prediction intervals (see Table 11). The prediction intervals of the SAC-SMA model are  
 1231 remarkably accurate and envelop about 44.3, 67.4 and 91.3% of the discharge observations for  
 1232  $\gamma = 50$ ,  $\gamma = 75$  and  $\gamma = 95\%$ , respectively.
- 1233 4. The forecast distributions of HYMOD and the SAC-SMA model do not only have a reasonable cov-  
 1234 erage but also are also remarkably reliable with tabulated  $R_l$  values of 0.895 and 0.942, respectively.  
 1235 These values are close to one (optimal) and suggest a well dispersed ensemble.

Table 11: Time-averaged values of the moments (skew and kurtosis) and performance metrics ( $R_l$ ,  $C_v$ ,  $C$  and  $W$ ) of the distribution forecasts simulated by HYMOD, Hmodel and the SAC-SMA model and RMSE,  $R^2$  and KGE scoring functions of the mean functional of the forecasts. The coverage  $C$  corresponds to significance levels  $\alpha = 0.5$ , 0.25 and 0.05 of Fig. 15 and  $d$  equals the number of unknown parameters.

Model	$d$	Skew	Kurt.	$R_l$	$C_v$	Coverage $C$ at $\gamma$			$W$	RMSE	$R^2$	KGE
						0.50	0.75	0.95				
HYMOD	5	1.485	10.60	0.105	0.796	0.385	0.615	0.855	2.401	1.590	0.695	0.679
Hmodel	7	0.003	2.854	0.558	0.365	0.214	0.446	0.778	1.549	1.555	0.708	0.611
SAC-SMA	13	1.042	6.016	0.058	0.739	0.443	0.674	0.913	2.386	1.375	0.772	0.708

1236 Altogether, the results demonstrate that the GLUE methodology coupled with a normal likelihood and  
 1237 logarithmic scoring rule can provide a remarkably proficient forecast ensemble. This offers a simple  
 1238 method for obtaining probabilistic estimates of the simulated state variables and discharge. Parameter  
 1239 calibration may further enhance the ensemble.

### 1240 6.2.5 Case Study VII: The Flow Duration Curve

1241 One obvious hydrologic application of the scoring rules is the flow duration curve (FDC). This signature  
 1242 catchment characteristic relates the exceedance probability of streamflow,  $\mathbb{P}(X > x)$ , to its magnitude,  $x$ ,  
 1243 and plays a critical role in (among others) flood frequency analysis, hydrologic model diagnostics, water  
 1244 quality management and the design of hydroelectric power plants (Sadegh *et al.*, 2016). The FDC is  
 1245 known as the survival function  $S_X(x)$  in statistics and the reliability function  $R_X(x)$  in engineering. We

1246 adopt the nomenclature of reliability

$$\begin{aligned}
 1247 \quad R_X(x) &= \mathbb{P}(X > x) = \int_x^{\infty} f_X(t) dt \\
 1248 \quad &= 1 - \int_{-\infty}^x f_X(t) dt \\
 1249 \quad &= 1 - F_X(x),
 \end{aligned} \tag{80}$$

1251 and reconfirm that the FDC is the complement of the streamflow CDF,  $F_X(x)$  (*Vogel and Fennessey*,  
 1252 1994). If we enter the above relationship,  $F_X(x) = 1 - R_X(x)$ , in the CRPS divergence function of  
 1253 Equation (67)

$$\begin{aligned}
 1254 \quad d_{\text{CRPS}}(P, Q) &= \int_0^{\infty} \left( (1 - R_P(z)) - (1 - R_Q(z)) \right)^2 dz, \\
 1255
 \end{aligned} \tag{81}$$

1256 we yield the divergence function of the supposed continuous ranked exceedance probability score

$$\begin{aligned}
 1257 \quad d_{\text{CREPS}}(P, Q) &= \int_0^{\infty} (R_Q(z) - R_P(z))^2 dz = d_{\text{CRPS}}(P, Q).
 \end{aligned} \tag{82}$$

We do not derive an expression for the actual scoring rule,  $S_{\text{CREPS}}(P, \omega)$ , of the CREPS divergence. This score will yield the same rankings of the simulated flow duration curves. Furthermore,  $d_{\text{CREPS}}(P, Q)$  is nonnegative and zero only when  $R_P = R_Q$ . We draw inspiration from *Thorarinsdottir et al.* (2013) and decompose the CREPS divergence into a term that summarizes the variability between the reliability functions of  $P$  and  $Q$ , and two other terms that measure the variability within the FDCs of  $P$  and  $Q$

$$\begin{aligned}
 1260 \quad d_{\text{CREPS}}(P, Q) &= \mathbb{E}_{P,Q}[|y - \omega|] - \frac{1}{2} (\mathbb{E}_P[|y - y^*|] - \mathbb{E}_Q[|\omega - \omega^*|]),
 \end{aligned} \tag{83}$$

1261 where  $(\omega, \omega^*)$  and  $(y, y^*)$  are independent copies of the measured and simulated discharge records,  
 1262 respectively. For a measured  $\omega_1, \dots, \omega_n$  and simulated  $y_1, \dots, y_n$  streamflow time series, we use a Monte  
 Carlo estimate of Equation (83)

$$\begin{aligned}
 1263 \quad d_{\text{CREPS}}(P, Q) &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n |y_i - \omega_j| - \frac{1}{2} \frac{1}{n^2} \left( \sum_{i=1}^n \sum_{j=1}^n |y_i - y_j| - \sum_{i=1}^n \sum_{j=1}^n |\omega_i - \omega_j| \right) \\
 1264 \quad &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n |y_i - \omega_j| - \frac{1}{2} \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \{ |y_i - y_j| - |\omega_i - \omega_j| \},
 \end{aligned} \tag{84}$$

1265 The above function is easy to compute and satisfies all the properties of a score divergence deemed  
 1266 desirable by *Ferson et al.* (2008). The  $d_{\text{CREPS}}(P, Q)$  function is (i) mathematically well behaved and  
 1267 understood;  $d_{\text{CREPS}}(P, Q) > 0$  unless  $R_P = R_Q$  then  $d_{\text{CREPS}}(P, Q) = 0$ , (ii) expressed in physical units  
 1268 (discharge, mm/d), (iii) sensitive to all moments of the FDC, not just mean and variance, and (iv) equal  
 1269 to the absolute error  $d_{\text{CREPS}}(\delta_y, \delta_\omega) = |\delta_y - \delta_\omega|$  between two point measures,  $\delta_y$  and  $\delta_\omega$ .

1271 We revisit the GLUE ensemble and compute Equation (84) for each of the  $m = 25,000$  simulated discharge  
 1272 records of the SAC-SMA model for the Leaf River basin. Figure 16 compares the measured FDC of the  
 1273 3,000-day training data record against simulated reliability functions of the fifty best ensemble members  
 1274 with lowest values of  $d_{\text{CREPS}}(P, Q)$  in Equation (84). To benchmark the CREPS divergence, we also present  
 1275 scatter diagrams of (b)  $d_{\text{CREPS}}(P, Q)$  and the mean taxicab distance,  $\bar{d}_T(P, Q) = \frac{1}{n} \sum_{t=1}^n |\omega'_t - y'_t(\boldsymbol{\theta}_i)|$ ,  
 1276 where  $\omega'_1, \dots, \omega'_n$  and  $y'_1, \dots, y'_n$  are ordered point measures of the measured and simulated discharge  
 1277 records, respectively, and (c)  $d_{\text{CREPS}}(P, Q)$  and the log-likelihood  $\ell(\boldsymbol{\theta}_i|\boldsymbol{\omega})$  in Equation (79) of the  
 1278 simulated discharge record. Each gray square corresponds to a different member of the GLUE ensemble.  
 The simulated FDCs of the SAC-SMA model are in close agreement with the actual reliability function

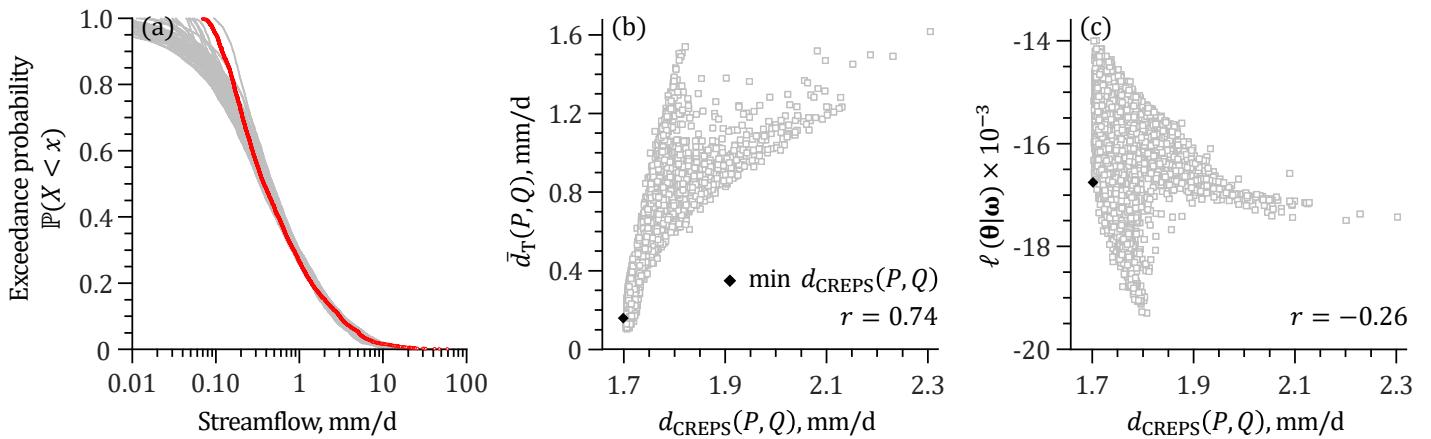


Figure 16: Preliminary results of the cumulative ranked exceedance probability divergence score: (a) measured FDC (red dots) and SAC-SMA model simulated reliability functions (gray lines) of the  $K = 50$  discharge records of the ensemble with lowest values of the CREPS divergence and (b,c) bivariate scatter plots of the CREPS divergence and (b) mean taxicab distance,  $\bar{d}_T(P, Q)$  and (c) normal log-likelihood  $\ell(\boldsymbol{\theta}|\boldsymbol{\omega})$  of Equation (79) using all discharge records of the GLUE ensemble. The black triangle corresponds to the minimum CREPS divergence.

1279  
 1280 of the Leaf River. The logarithmic scale emphasizes the discrepancies for streamflows smaller than about  
 1281 0.1 mm/d. This apparent mismatch will disappear with a linear streamflow scale. Alternatively, we  
 1282 can apply a nonnegative weight function  $\int_0^\infty w(z) dz < \infty$  to the squared residuals of the measured and  
 1283 simulated FDCs

$$d_{\text{WCRES}}(P, Q) = \int_0^\infty (R_Q(z) - R_P(z))^2 w(z) dz, \quad (85)$$

1284 and emphasize description of the relative frequencies of the lowest streamflows. When  $P, Q \in \mathcal{P}_1$  the  
 1285 above expression is a valid generalization of the CREPS divergence score (Matheson and Winkler, 1976).  
 1286 The CREPS divergence of the measured and simulated reliability functions correlates reasonably well  
 1287 with the taxicab distance of  $R_Q$  and  $R_P$ . But the CREPS divergence is not equivalent to a Manhattan

(and Euclidean) distance, and as a result, we find minimum values of  $d_{\text{CRPS}}(P, Q)$  over a range of taxicab distances. The CREPS divergence correlates poorly with the log-likelihood of the SAC-SMA simulated discharge records and attains a minimum value of about 1.7 mm/d well removed from the likelihood maximum of  $\ell(\boldsymbol{\theta}|\boldsymbol{\theta}) \approx -14 \times 10^3$ . This demonstrates once again that purely statistical metrics of the goodness of fit may compromise the ability of the model to reproduce hydrologically relevant signatures of catchment functioning. This testifies to the added value of process-based diagnostic measures embodied in a *strictly proper* score divergence.

### 6.3 Interval scoring rules

If so desired one can summarize probabilistic forecasts of a continuous quantity using predictive quantiles. Suppose the forecaster quotes quantiles  $\mathbf{r} = (r_1, \dots, r_z)^\top$  and  $x$  materializes, then he or she will be rewarded by the score  $S(\mathbf{r}; x)$ . Then, the expected score  $\mathcal{S}(\mathbf{r}; P)$  under probability measure  $P \in \mathcal{P}$  is equal to (*Gneiting and Raftery, 2007*)

$$\mathcal{S}(\mathbf{r}; P) = \int S(\mathbf{r}; x) dP(x). \quad (86)$$

If  $q_1, \dots, q_z$  are the true quantiles for the class  $\mathcal{P}$  of Borel probability measures on  $\mathbb{R}$  then a scoring rule is proper if (*Cervera and Muñoz, 1996*)

$$\mathcal{S}(q_1, \dots, q_z; P) \geq \mathcal{S}(r_1, \dots, r_z; P) \quad (87)$$

for all real numbers  $r_1, \dots, r_z$  and  $P \in \mathcal{P}$ . *Gneiting and Raftery* (2007) present a general form of a scoring rule for quantiles. We focus our attention on the coverage  $C$  in Table 1 and formulate this insufficient performance metric as a *proper* scoring rule of the  $100(1 - \alpha)\%$  prediction interval

$$S_{\text{IS}}^\alpha(P, \omega) = (l - u) - \frac{2}{\alpha}(l - \omega)\mathbb{1}\{\omega \leq l\} - \frac{2}{\alpha}(\omega - u)\mathbb{1}\{\omega \geq u\}, \quad (88)$$

where  $l = F_P^{-1}(\alpha/2)$  and  $u = F_P^{-1}(1 - \alpha/2)$  denote the lower and upper endpoints of the predictive quantiles at significance levels  $\alpha/2$  and  $1 - \alpha/2$ , respectively. The interval score  $S_{\text{IS}}^\alpha(P, \omega)$  is positively oriented and incurs a penalty, the size of which depends on the significance level  $\alpha$ , if the observation  $\omega$  materializes outside the  $[u, l]$  prediction interval. From the first term on the right-hand side of Equation (88) it is evident that  $S_{\text{IS}}^\alpha(P, \omega)$  rewards narrow prediction intervals.

We perform a simple numerical experiment to demonstrate that  $S_{\text{IS}}^\alpha(P, \omega)$  is a *proper* scoring rule. We draw at random  $n = 10^4$  observations,  $\omega_1, \dots, \omega_n$ , from a standard normal distribution. From tabulated

critical values, we expect that about 95% of these observations lie in the interval  $-1.96 < \omega < 1.96$ . Thus, if we fix  $\alpha = 0.05$  in Equation (88) then  $l = -1.96$  and  $u = 1.96$  should maximize the expected value of  $S_{\text{IS}}^{\alpha}(P, \omega)$ . We verify this assertion in Figure 17 and plot the mean value of the interval score as function of the lower endpoint  $l \in [-5, -0.1]$  using  $u = 1.96$  and  $\alpha = 0.05$  (red),  $\alpha = 0.25$  (blue) and  $\alpha = 0.50$  (green). The colored lines display a strong dependency of the mean interval score  $S_{\text{IS}}^{\alpha}(P, \omega)$  on the choice

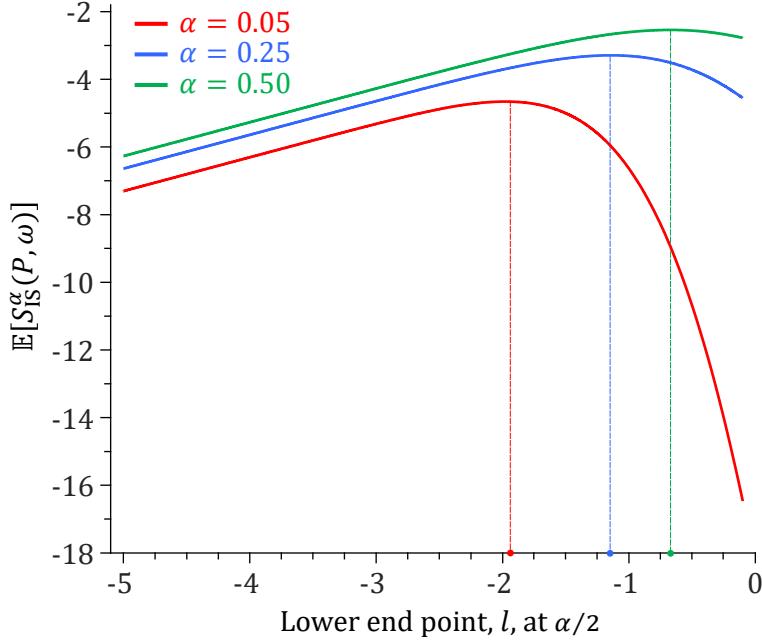


Figure 17: Traces of the expected value of the interval score  $S_{\text{IS}}^{\alpha}(P, \omega)$  as function of the lower endpoint  $l$  of the  $100(1 - \alpha)\%$  prediction interval using  $\alpha = 0.05$  (red),  $\alpha = 0.25$  (blue) and  $\alpha = 0.50$  (green). The colored dots are a projection of the maximum interval score on the  $x$ -axis.

of the lower endpoint  $l = F_P^{-1}(\alpha/2)$ . The interval score achieves its largest value, on average, when the forecaster quotes the *true* lower endpoints,  $l = -1.96$ ,  $l = -1.15$  and  $l = -0.67$  of the  $100(1 - \alpha)\%$  prediction intervals of  $\omega \sim \mathcal{N}(0, 1)$  at significance levels  $\alpha = 0.05$ ,  $\alpha = 0.25$  and  $\alpha = 0.50$ , respectively. This encourages the forecaster to be honest and volunteer his or her true beliefs. To guarantee an accurate description of the *true* forecast distribution one would need to compute  $S_{\text{IS}}^{\alpha}(\mathbf{P}, \omega)$  at many different  $\alpha \in (0, 1)$  values. But as in Christoffersen (1998) this is a daunting task, most certainly in the presence of trade-offs in the interval score of different prediction intervals. When confronted with a time series of forecasts  $P_1, \dots, P_n$  we work with the time-averaged interval score

$$\bar{S}_{\text{IS}}^{\alpha}(\mathbf{P}, \boldsymbol{\omega}) = \frac{1}{n} \sum_{t=1}^n S_{\text{IS}}^{\alpha}(P_t, \omega_t). \quad (89)$$

## 1330 6.4 Multivariate Forecasts

1331 Up until now, we have considered forecast distributions of only a single variable of interest, say, discharge,  
 1332 and verifying observations measured at different times. The overall skill score,  $\bar{S}(\mathbf{P}, \boldsymbol{\omega})$ , is then a  
 1333 time-averaged mean score as defined in Equation (74). We can expand this approach to multi-variable  
 1334 forecasts using marginalization by treating the distribution of each variable, say discharge, soil moisture  
 1335 content and respective aspects of stream water chemistry, separately. But for such multi-variable forecasts,  
 1336  $P \in \mathcal{P} \in \mathbb{R}^b$ , we can also resort to multivariate scoring rules such as the energy score of *Gneiting and*  
 1337 *Raftery* (2007)

$$1338 S_{\text{ES}}(P, \boldsymbol{\omega}) = \frac{1}{2} \mathbb{E}_P [\|\mathbf{y} - \mathbf{y}^*\|_2^\eta] - \mathbb{E}_P [\|\mathbf{y} - \boldsymbol{\omega}\|_2^\eta], \quad (90)$$

1339 where  $\|\cdot\|_2$  is the Euclidean norm on  $\mathbb{R}^b$ ,  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_b)^\top$ ,  $\eta \in (0, 2)$ , and  $\mathbf{y} = (y_1, \dots, y_b)^\top$  and  
 1340  $\mathbf{y}^* = (y_1^*, \dots, y_b^*)^\top$  are independent copies of the  $b$ -variate distribution,  $P \in \mathcal{P}_\eta$ . For  $b = 1$ , the above  
 1341 expression reduces to  $S_{\text{ES}}(P, \omega)$  in Equation (68). The multivariate form of the energy score is a *strictly*  
 1342 *proper* score (*Székely*, 2003) and reduces to

$$1343 S_{\text{SE}}(P, \boldsymbol{\omega}) = -\|\boldsymbol{\mu}_P - \boldsymbol{\omega}\|_2^2, \quad (91)$$

1344 for  $\eta \rightarrow 2$ , where  $\boldsymbol{\mu}_P = (\mu_{1P}, \dots, \mu_{bP})^\top$  signifies the mean of the forecast distribution. The squared error  
 1345 can also be written as a vector-inner product,  $S_{\text{SE}}(P, \boldsymbol{\omega}) = -(\boldsymbol{\mu}_P - \boldsymbol{\omega})^\top(\boldsymbol{\mu}_P - \boldsymbol{\omega})$ . Analogous to the  
 1346 numerical definition of the CRPS in Equation (64), the ES may be approximated using a large collection  
 1347  $\{\mathbf{y}_1, \dots, \mathbf{y}_m\}$  of  $m$  samples of the forecast distribution  $P$  (*Grimit et al.*, 2006)

$$1348 S_{\text{ES}}(P, \boldsymbol{\omega}) = \frac{1}{2m^2} \sum_{i=1}^m \sum_{j=1}^m \|\mathbf{y}_i - \mathbf{y}_j\|_2^\eta - \frac{1}{m} \sum_{i=1}^m \|\mathbf{y}_i - \boldsymbol{\omega}\|_2^\eta. \quad (92)$$

1349 *Dawid* (1998) and *Dawid and Sebastiani* (1999) studied scoring rules that depend only on the mean,  
 1350  $\boldsymbol{\mu}_P \in \mathbb{R}^{b \times 1}$ , and the dispersion or covariance matrix,  $\boldsymbol{\Sigma}_P \in \mathbb{R}^{b \times b}$  of the forecast distribution  $P$ . Their  
 1351 divergence function,  $d_{\text{DDS}}(P, Q)$ , in Equation (5) is linked to the *proper* scoring rule (*Dawid and Sebastiani*,  
 1352 1999)

$$1353 S_{\text{DSS}}(P, \boldsymbol{\omega}) = -\log_e(|\boldsymbol{\Sigma}_P|) - (\boldsymbol{\omega} - \boldsymbol{\mu}_P)^\top \boldsymbol{\Sigma}_P^{-1} (\boldsymbol{\omega} - \boldsymbol{\mu}_P), \quad (93)$$

1354 and generalized entropy function

$$1355 H(P) = -\log_e(|\boldsymbol{\Sigma}_P|) - b. \quad (94)$$

1356 Note that the DSS is simply equal to the unnormalized log-likelihood of a multivariate normal density.  
 1357 For a univariate forecast,  $b = 1$ , the DSS simplifies to  $S_{\text{DSS}}(P, \omega) = -\log_e(\sigma_P^2) - (\omega - \mu_P)^2/\sigma_P^2$ , where,  
 1358 again,  $\mu_P$  and  $\sigma_P^2$  are the mean and variance of the forecast distribution.

A multivariate forecast does not necessarily have to involve multiple different variables. It can also consist of a single variable but forecasted at many different sites. Suppose  $i, j \in (1, \dots, b)$  are single indices for the points in a square or rectangular grid made up of  $b$  sites. Scheuerer and Hamill (2015) investigate the accuracy of the forecast distribution  $P$  in describing the spatial structure of gridded measurements  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_b)^\top$  of 80-meter wind speed forecasts. They introduced the so-called variogram score of order  $\varsigma > 0$

$$S_{VS}^\varsigma(P, \boldsymbol{\omega}) = - \sum_{i=1}^b \sum_{j=1}^b w_{ij} (|\omega_i - \omega_j|^\varsigma - \mathbb{E}_P [|y_i - y_j|^\varsigma])^2, \quad (95)$$

where  $w_{ij} \geq 0$  is a nonnegative weight attached to the  $(i, j)^{\text{th}}$  pair of sites and  $y_i$  and  $y_j$  are the  $i^{\text{th}}$  and the  $j^{\text{th}}$  elements (sites) of a random vector that is distributed according to  $P$ . The weights allows one to emphasize or downplay specific aspects of the distribution forecast. The variogram scoring rule is analogous to the squared error  $S_{SE}(P, \boldsymbol{\omega})$  of Equation (91) using residuals of the powered differences of pairs of measurement sites and pairs of forecast sites. For  $\varsigma = 2$  the powered difference is known as the semi-variance. When the forecast distribution  $P$  is given in the form of a  $m$ -member ensemble  $\{\mathbf{y}_1, \dots, \mathbf{y}_m\}$  the second term of  $S_{VS}^\varsigma(P, \boldsymbol{\omega})$  can be approximated by

$$\mathbb{E}_P [|y_i - y_j|^\varsigma] = \frac{1}{m} \sum_{k=1}^m |y_{ik} - y_{jk}|^\varsigma. \quad (96)$$

The variogram scoring rule is more sensitive than existing scoring rules to multivariate site dependencies. But as the variogram score admits site differences, it is insensitive to the location of the distribution forecast, thus, is a *proper* and not *strictly proper* scoring rule. Forecast distributions that differ only in their mean will receive the same values of  $S_{VS}^\varsigma(P, \boldsymbol{\omega})$ .

We could conveniently discard time as controlling variable of the rainfall-runoff transformation and treat the discharge time series of HYMOD, Hmodel and the SAC-SMA model as a multivariate forecast  $\mathbf{y} = (y_1, \dots, y_b)^\top$  with dimensionality  $b$  equal to the length  $n = 3,000$  of the simulated records. The  $m$ -member ensemble  $\{\mathbf{y}_1, \dots, \mathbf{y}_m\}$  will produce values of  $S_{ES}(P, \boldsymbol{\omega})$ ,  $S_{DSS}(P, \boldsymbol{\omega})$  and  $S_{VS}^\varsigma(P, \boldsymbol{\omega})$  according to Equations (92), (94) and (95), respectively.

## 7 Decomposition of scoring rules for categorical forecasts

Strictly proper scoring rules condense the accuracy of the distribution forecast  $P$  to a single reward oriented value while retaining attractive statistical properties. This compression to a single reward value

1386 simplifies tasks such as forecast verification, model evaluation and likelihood function selection (e.g.  
 1387 *Vrugt et al.* 2022), but makes it difficult to pinpoint exactly which attributes of the distribution forecast  
 1388 are deficient and in need of improvement. As a result, *strictly proper* scoring rules lack diagnostic power  
 1389 on how to improve the overall consistency, accuracy and precision of model forecasts. Thus, it may be  
 1390 beneficial to decompose scoring rules into their constituent parts. Decomposition of the expected loss  
 1391 into a calibration and refinement loss has facilitated the development of calibration methods (*Bella et al.*,  
 1392 2013). The refinement loss consists of an uncertainty term and a resolution term (*DeGroot and Fienberg*,  
 1393 1983; *Murphy*, 1973). *Kull and Flach* (2015) presents a decomposition of the logarithmic and Brier  
 1394 scoring rules into an epistemic and irreducible (aleatoric) loss term. Next, we review the decomposition  
 1395 of *strictly proper* scoring rules into an uncertainty, resolution and reliability term. These components  
 1396 relate directly to forecast attributes that are deemed desirable on grounds independent of the scoring  
 1397 rules themselves and provide an epistemological justification of measuring forecast quality by *strictly*  
 1398 *proper* scoring rules (*Bröcker*, 2009).

## 1399 7.1 Theory

1400 Let  $\Omega = \{1, 0\}$  be the sample space of a binary event of *rain* or *no rain*. Let the quoted probability  
 1401  $p = p(\mathcal{D})$  of *rain* be a function of the data  $\mathcal{D}$  available to the forecaster up to a certain lead time, where  
 1402  $p \in [0, 1]$ . A forecasting scheme models the relationship between the data  $\mathcal{D}$  and the quantity of interest  
 1403 or forecast (see *Murphy and Winkler* 1987; *Murphy* 1993, 1996 for a related discussion). Once we observe  
 1404  $\omega \in \Omega$ , we can assign a score  $S(p, \omega) : \mathcal{P}_2 \times \Omega \rightarrow \mathbb{R}$  to the prediction. Thus,  $\omega$  takes on values of 0 (*no*  
 1405 *rain*) and 1 (*rain*). The law of total expectation states that if  $S(p, \omega)$  and  $\mathcal{D}$  are random variables on the  
 1406 same probability space then

$$1407 \quad \mathbb{E}[S(p, \omega)] = \mathbb{E}[\mathbb{E}[S(p, \omega)|\mathcal{D}]]. \quad (97)$$

1408 Suppose we use the Brier or quadratic score,  $S_{\text{QS}}(p, \omega) = -\sum_{k=1}^2 (\delta_{\omega k} - p_k)^2$ , where  $\delta_{\omega k} = 1$  if  $\omega = k$  and  
 1409  $\delta_{\omega k} = 0$  otherwise. Then we can decompose the conditional expectation

$$\begin{aligned}
 1410 \quad \mathbb{E}[S_{\text{QS}}(p, \omega)|\mathcal{D}] &= \mathbb{E}[-(\omega - p(\mathcal{D}))^2|\mathcal{D}] \\
 1411 &= -\mathbb{E}[(\omega - \mathbb{E}[\omega|\mathcal{D}] + \mathbb{E}[\omega|\mathcal{D}] - p(\mathcal{D}))^2|\mathcal{D}] \\
 1412 &= -\text{Var}[\omega|\mathcal{D}] - (\mathbb{E}[\omega|\mathcal{D}] - p(\mathcal{D}))^2,
 \end{aligned} \quad (98)$$

1414 and insert Equation (98) into Equation (97) to yield

$$\begin{aligned}\mathbb{E}[S_{QS}(p, \omega)] &= \mathbb{E}[-\text{Var}[\omega|\mathcal{D}] - (\mathbb{E}[\omega|\mathcal{D}] - p(\mathcal{D}))^2] \\ &= -\mathbb{E}[\text{Var}[\omega|\mathcal{D}]] - \mathbb{E}[(\mathbb{E}[\omega|\mathcal{D}] - p(\mathcal{D}))^2] \\ &= -\text{Var}[\omega] + \text{Var}[\mathbb{E}[\omega|\mathcal{D}]] - \mathbb{E}[(p(\mathcal{D}) - \mathbb{E}[\omega|\mathcal{D}])^2],\end{aligned}\quad (99)$$

1419 where  $\mathbb{E}[\omega|\mathcal{D}]$  is simply equal to the conditional probability of *rain* and  $p(\mathcal{D})$  equals the unconditional  
1420 *rain* probability. Bröcker (2009) generalized the above decomposition to a generic *strictly proper* scoring  
1421 rule,  $S(\mathbf{p}, \omega) : \mathcal{P}_m \times \Omega \rightarrow \mathbb{R}$  of a categorical forecast of  $m \geq 2$  events to yield

$$\mathbb{E}[S(\mathbf{p}, \omega)] = \underbrace{H(\bar{\mathbf{p}})}_1 + \underbrace{\mathbb{E}[d(\bar{\mathbf{p}}, \boldsymbol{\pi})]}_2 - \underbrace{\mathbb{E}[d(\mathbf{p}, \boldsymbol{\pi})]}_3, \quad (100)$$

1422 where  $\bar{\mathbf{p}} = (\bar{p}_1, \dots, \bar{p}_m)^\top$  is the unconditional probability of  $\omega$  (so-called climatology in meteorology) and  
1423  $\pi_k = \mathbb{P}(\omega = k|\mathbf{p})$  signifies the conditional probability of observation  $\omega$  for the probabilities  $\mathbf{p}$  quoted,  
1424  $k = (1, \dots, m)$ . Hence,  $\boldsymbol{\pi}$  is a mapping ( $m \times m$  matrix) which specifies for every  $\omega \in \Omega$  a probability  
1425 measure on  $\mathcal{P}_m$ . The three terms are nonnegative for *strictly proper* scoring rules and referred to as (1)  
1426 uncertainty (of  $\omega$ ), (2) resolution (or sharpness) and (3) reliability (Bröcker, 2009). The entropy and  
1427 resolution have a positive effect on  $\mathbb{E}[S(\mathbf{p}, \omega)]$ , whereas the reliability decreases the expected value of  
1428 the score. Note that the minus sign of  $-\text{Var}[\omega]$  in Equation (99) has disappeared from the uncertainty  
1429 term as Bröcker (2009) uses  $-H(\bar{\mathbf{p}})$  for the entropy function. Furthermore, the use of positively oriented  
1430 scoring rules reverses the signs in front of the resolution and reliability terms.

1433 The first term of the decomposition is uncertainty and quantifies our state of knowledge in the absence of  
1434 an underlying theory to generate forecasts. This is equivalent to the expected score of the average event  
1435 frequencies or climatology as forecast probabilities. This term is independent of the forecasting system  
1436 and depends only on the statistics of the observations (Christensen, 2015). The second term,  $\mathbb{E}[d(\bar{\mathbf{p}}, \boldsymbol{\pi})]$   
1437 or resolution measures the mean divergence of the conditional event probabilities  $\boldsymbol{\pi}$  from their average  
1438 probabilities  $\bar{\mathbf{p}}$  and, thus is a proxy for the variance of  $\boldsymbol{\pi}$  over the sample  $\Omega$  or data  $\mathcal{D}$  space. At zero  
1439 resolution  $\boldsymbol{\pi}$  is always equal to the climatology (or prior mean)  $\bar{\mathbf{p}}$  and, thus, the data  $\mathcal{D}$  provides no useful  
1440 information. Thus, in accordance with the sharpness principle of (Gneiting and Raftery, 2007) larger  
1441 values of the resolution are preferred and reflect case-dependent probabilistic forecasts. The third and  
1442 last term, reliability, measures the average deviation of the probabilistic forecast  $\mathbf{p}$  from the conditional  
1443 event probabilities  $\boldsymbol{\pi}$ . This is a measure of the statistical consistency of a forecast and evaluates whether

1444 the quoted forecast probabilities are in agreement with the materialized event frequencies. Thus, the  
 1445 reliability penalizes poorly calibrated forecasts as  $\mathbf{p}$  will not match  $\boldsymbol{\pi}$ .

1446 Equation (100) is a generalisation of the well known decomposition of the Brier score of *Murphy* (1973)

1447 1. uncertainty =  $\bar{p}(1 - \bar{p})$       2. resolution =  $\mathbb{E}[(\bar{p} - \pi_1)^2]$       3. reliability =  $\mathbb{E}[(p - \pi_1)^2]$ ,      (101)

1448 where  $\pi_1 = \pi$  is the conditional probability of *rain* given  $p$ ,  $\pi_1 = \mathbb{P}(x = 1|p)$ . We refer to *Weijns et al.*  
 1449 (2010b) for a reliability-resolution-uncertainty decomposition of the KL-divergence, the divergence of the  
 1450 logarithmic scoring rule.

## 1451 7.2 Case Study VII: Discharge Forecast Ensemble

1452 We illustrate the analytic decomposition of Equation (100) by application to the multi-model ensemble  
 1453 of discharge forecasts displayed in Figure 11. In doing so we must first convert discharge to a categorical  
 1454 variable with number of possible outcomes  $m$  equal to the ensemble size,  $K = 8$ . In this discrete sample  
 1455 space,  $\Omega = \{1, \dots, m\}$ , the measured daily discharge  $\omega_t$  at  $t \geq 1$  coincides with the "best" discharge  
 1456 forecast among the watershed models. The index  $\hbar$  of this "best" forecast at each time  $t$  is the index of  
 1457 the minimum entry of the  $m$ -vector of absolute residuals

1458 
$$\hbar = \arg \min_{\hbar \in \{1, \dots, K\}} |\omega_t - y_{\hbar t}|, \quad (102)$$

1459 between the measured and forecasted discharges. Suppose the vector  $(\hbar_1, \dots, \hbar_n)^\top$  stores the indices of  
 1460 the best models for our training record of  $n = 3,000$  days, then the event frequencies  $\bar{\mathbf{p}} = (\bar{p}_1, \dots, \bar{p}_m)$   
 1461 may be computed using

1462 
$$\bar{p}_k = \mathbb{P}(\omega = y_k) = \frac{1}{n} \sum_{t=1}^n \mathbb{1}\{\hbar_t = k\}, \quad (103)$$

1463 and the conditional forecast probabilities equal

1464 
$$\begin{aligned} \pi_{jk} &= \mathbb{P}(\omega = y_j | y_k) = \frac{\sum_{t=2}^n \mathbb{1}\{\hbar_t = j | \hbar_{t-1} = k\}}{\sum_{t=1}^n \mathbb{1}\{\hbar_t = k\}} \\ &= \frac{\sum_{t=2}^n \mathbb{1}\{\hbar_t = j | \hbar_{t-1} = k\}}{n \bar{p}_k}, \end{aligned} \quad (104)$$

1467 where  $j, k = (1, \dots, K)$ . Table 12 reports the unconditional and conditional probabilities of the watershed  
 1468 models.

Table 12: Unconditional,  $\bar{\mathbf{p}}$ , and conditional,  $\boldsymbol{\pi}$ , probabilities of the watershed models estimated from the 3,000-day training data record:  $\pi_{jk} = \mathbb{P}(\omega = y_j | y_k)$  is the probability of  $y_j$  given that  $y_k$  is the best forecast in the ensemble at the previous time.

Model	ABC	GR4J	HYMOD	TOPMO	AWBM	NAM	HBV	SAC-SMA
$\bar{\mathbf{p}}$	0.064	0.148	0.088	0.101	0.080	0.142	0.175	0.203
$\boldsymbol{\pi}$	ABC	<b>0.267</b>	0.050	0.075	0.056	0.046	0.056	0.048
	GR4J	0.157	<b>0.534</b>	0.098	0.076	0.087	0.033	0.061
	HYMOD	0.099	0.063	<b>0.343</b>	0.102	0.083	0.035	0.057
	TOPMO	0.099	0.070	0.094	<b>0.383</b>	0.054	0.059	0.069
	AWBM	0.037	0.043	0.053	0.063	<b>0.442</b>	0.052	0.033
	NAM	0.099	0.045	0.087	0.106	0.062	<b>0.609</b>	0.038
	HBV	0.105	0.086	0.094	0.102	0.033	0.061	<b>0.608</b>
	SAC-SMA	0.136	0.110	0.155	0.112	0.192	0.094	<b>0.531</b>
$\sum_{j=1}^8 \pi_{jk}$	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000

<sup>1469</sup> The unconditional forecast probabilities tend to increase with model complexity and is largest for the  
<sup>1470</sup> SAC-SMA model. The conditional forecast probabilities are largest on the main diagonal (in bold)  
<sup>1471</sup> confirming that a model's probability is largest conditional on it having the best forecast in the ensemble.  
<sup>1472</sup> Note that each column of  $\boldsymbol{\pi}$  sums to unity.

<sup>1473</sup> Table 13 presents the expected value of the quadratic score, entropy, resolution and reliability of the  
<sup>1474</sup> BMA mixture distribution for the PDFs of Table 8 with a constant group and sole forecast variance,  
<sup>1475</sup> respectively. In each case, we set the forecast probabilities  $\mathbf{p} = (p_1, \dots, p_K)^\top$  equal to the maximum  
<sup>1476</sup> likelihood values of the weights  $\beta_1, \dots, \beta_K$  of the BMA mixture distribution.

Table 13: Time-averaged expected value of the *strictly proper* quadratic scoring rule (see Table 6), entropy, resolution and reliability for the BMA density forecasts of Equation (70) using the normal, lognormal, generalized normal, gamma and Weibull distributions with a constant group and sole variance, respectively. The bottom row completes the decomposition of Equation (100).

	Normal <sup>a</sup>		Lognormal		Gen. Normal		Gamma		Weibull	
	group	sole	group	sole	group	sole	group	sole	group	sole
$\mathbb{E}[S_{QS}(\mathbf{p}, \omega)]$	0.027	0.100	0.031	0.033	0.108	0.108	0.073	0.116	0.066	0.071
$H(\bar{\mathbf{p}})$	0.142	0.142	0.142	0.142	0.142	0.142	0.142	0.142	0.142	0.142
$\mathbb{E}[d(\bar{\mathbf{p}}, \boldsymbol{\pi})]$	0.162	0.162	0.162	0.162	0.162	0.162	0.162	0.162	0.162	0.162
$\mathbb{E}[d(\mathbf{p}, \boldsymbol{\pi})]$	0.277	0.204	0.272	0.271	0.196	0.196	0.231	0.188	0.238	0.233
Sum	0.027	0.100	0.031	0.033	0.108	0.108	0.073	0.116	0.066	0.071

<sup>a</sup> We fix  $\tau = 2$  in the generalized normal density of Table 8

<sup>1477</sup> The tabulated data confirm the decomposition of Equation (100). The expected value of the quadratic

scoring rule (top row) is indeed equal to the entropy,  $H(\bar{\mathbf{p}})$ , plus the resolution,  $\mathbb{E}[d(\bar{\mathbf{p}}, \boldsymbol{\pi})]$ , minus the reliability,  $\mathbb{E}[d(\mathbf{p}, \boldsymbol{\pi})]$ . The first two terms of this decomposition depend only on the unconditional and conditional event frequencies,  $\bar{\mathbf{p}}$  and  $\boldsymbol{\pi}$ , respectively, hence, do not differ among the component functions of the BMA forecast distribution. With exception of the generalized normal density, the use of a model-dependent forecast variance increases forecast reliability. The BMA model probabilities  $\mathbf{p} = (p_1, \dots, p_K)^\top$  derived from the gamma distribution maximize the expected value of the quadratic scoring rule. This conclusion contradicts earlier findings reported in Table 9 which linked the gamma distribution to the lowest value of the quadratic score. Thus, the use of the maximum likelihood solution of the gamma distribution leads to diametrically opposed conclusions about model adequacy. This testifies to the premise of this paper that only a complete evaluation of the forecast density provides an honest and fair assessment of model adequacy.

## 8 Outlook

Scoring rules lie at the heart of statistical theory and practice and guarantee a more robust and complete evaluation of hydrologic models. But scoring rules may satisfy other purposes as well. We briefly describe a few avenues for future work. For obvious reasons, this discussion is succinct and incomplete.

### 8.0.1 Localized Scoring Rules: Extreme Events

One may only be interested in evaluating certain aspects of the forecast distribution, for example, the probability of extreme events in the lower and/or upper tail of the distribution. The CREPS divergence  $d_{WCREPS}(P, Q)$  in Equation (85) allows differential weighting of the flow levels in the description of the flow duration curve. If all weights are strictly positive, then this localization will not sacrifice the strict propriety of the CREPS divergence. The nonnegative weight function  $w_{ij} = \max(0, 1 - \frac{1}{9}|i - j|^2)$  used by *Scheuerer and Hamill* (2015) in the application of the  $S_{VS}^c(P, \boldsymbol{\omega})$  favors an accurate probabilistic description of the powered differences between nearby sites over such differences from distant sites. But as all sites with  $|i - j| \geq 3$  attain a zero weight, the above formulation of the variogram score is locally proper at best. No localized scoring rule can differentiate between predictive densities on sets with a zero weight, thus, satisfy propriety only for the domain of interest, say extreme events (*Diks et al.*, 2011;

1504 *de Punder et al.*, 2022). A simple work-around is to assign a minimum nominal weight to each set. The  
1505 weight function also enables users to incorporate soft information and/or personal judgment in model  
1506 evaluation.

1507 **8.0.2 Flood frequency analysis**

1508 Flood frequency analysis usually involves the fitting of the parameters of some known probability  
1509 distribution to a training data record of log-transformed annual maxima discharges,  $\omega_1, \dots, \omega_n$ . The  
1510 marginal likelihood (Bayes factors) of the estimated model parameters will convey which assumption  
1511 about the mean of the distribution (constant or time-dependent) receives most support from the measured  
1512 data record (*Luke et al.*, 2017). This is one of many examples which may benefit from the application of  
1513 scoring rules for hypothesis testing of stationary and nonstationary flood frequency models. Suppose we  
1514 choose the Pearson type III distribution  $\mathcal{P}_{\text{III}}(\mu, \sigma^2, \rho)$  with mean  $\mu$ , variance  $\sigma^2$  and skewness  $\rho$  for the  
1515 log-transformed annual maxima discharges. If we reparameterize the location, shape and scale parameters  
1516 of the PIII distribution to  $\xi = \mu - 2\sigma/\rho$ ,  $a = 4/\rho^2$  and  $b = \frac{1}{2}\sigma|\rho|$ , respectively, then the PDF of  $P$   
1517 simplifies to (*Hosking and Wallis*, 1997; *Tegos et al.*, 2022)

$$f_P(x, \xi, a, b) = \frac{|x - \xi|^{a-1}}{b^a \Gamma(a)} \exp(-b^{-1}|x - \xi|), \quad (105)$$

1519 where  $x, \xi, a, b \in \mathbb{R}$ ,  $a > 0$  and  $b > 0$ . If  $\rho > 0$  then  $x \in (\xi, \infty)$ , otherwise for  $\rho < 0$  we yield  
1520  $x \in (-\infty, \xi)$ . The logarithmic score,  $S_{\text{LS}}(P, \omega) = \log_b(f_P(\omega))$ , of distribution forecast  $P = \mathcal{P}_{\text{III}}(\xi, a, b)$   
1521 for the materialized outcome  $\omega_t$  is now equal to

$$S_{\text{LS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega_t) = (a - 1) \log_e(|\omega_t - \xi|) - a \log_e(b) - \log_e(\Gamma(a)) - b^{-1}|\omega_t - \xi|, \quad (106)$$

1523 with Euler's number as logical choice for the base of the score. The model that achieves the highest value  
1524 of the *strictly proper* logarithmic score is preferred by the data record. The QS, PSS and SS are readily  
1525 computed as well but an analytic expression for the CRPS of a PIII distribution forecast  $P = \mathcal{P}_{\text{III}}(\xi, a, b)$   
1526 is more involved (see Appendix I) and results in

$$\begin{aligned} S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega) &= 2 \times \frac{4^{-a}b}{B(a, a)} - ab + |\omega - \xi| + 2abF_{\mathcal{G}}(|\omega - \xi|, a + 1, b) \\ &\quad - 2|\omega - \xi|F_{\mathcal{G}}(|\omega - \xi|, a, b), \end{aligned} \quad (107)$$

1530 where  $B(u, v) = \Gamma(u)\Gamma(v)/\Gamma(u + v)$  is the beta function of the first kind and  $F_{\mathcal{G}}(z, a, b)$  is the CDF of  
1531 the gamma distribution  $\mathcal{G}(a, b)$  with shape and scale parameters,  $a > 0$  and  $b > 0$ , respectively. The

1532 first term in the CRPS of  $P = \mathcal{P}_{\text{III}}(\xi, a, b)$  is equal to the first term,  $\frac{1}{2}\mathbb{E}_P[|y - y^*|]$ , of Equation (63) and  
 1533 related to the well-known concentration index  $G$  of Corrado Gini for measuring the degree of inequality  
 1534 in the distribution of income and wealth (*McDonald and Jensen*, 1979; *Scheuerer and Möller*, 2015)

$$1535 \quad \frac{1}{2}\mathbb{E}_P[|y - y^*|] = 2 \times \frac{4^{-a}b}{B(a, a)} = abG = ab \frac{\Gamma(a + \frac{1}{2})}{\sqrt{\pi}\Gamma(a + 1)}. \quad (108)$$

1536 As  $\Gamma(\frac{1}{2}) = \sqrt{\pi}$ , we yield the numerically more robust  $\frac{ab}{\pi}B(a + \frac{1}{2}, \frac{1}{2})$  for the first term of Equation (107).  
 1537 For completeness, Appendix J presents an analytic form of the CRPS for a generalized extreme value  
 1538 distribution forecast  $P = \mathcal{GEV}(\mu, \sigma^2, \xi)$  (*Friederichs and Thorarinsdottir*, 2012).

### 1539 8.0.3 Bayesian model selection

1540 We revisit Bayes theorem in Equation (15) and look in more detail at the normalization constant  $p(\mathcal{D}|\mathcal{H})$

$$1541 \quad p(\mathcal{D}|\mathcal{H}) = \int_{\Theta} p(\mathcal{D}|\boldsymbol{\theta}, \mathcal{H})p(\boldsymbol{\theta}|\mathcal{H}) d\boldsymbol{\theta}. \quad (109)$$

1542 where  $L(\boldsymbol{\theta}|\mathcal{D}, \mathcal{H}) \equiv p(\mathcal{D}|\boldsymbol{\theta}, \mathcal{H})$  is the likelihood under hypothesis  $\mathcal{H}$  and  $p(\boldsymbol{\theta}|\mathcal{H})$  denotes the prior  
 1543 parameter distribution. The normalizing constant is also referred to as the integrated or marginal  
 1544 likelihood. Now for two competing hypothesis,  $\mathcal{H}_0$  and  $\mathcal{H}_1$ , of a finite or countable class of statistical  
 1545 models with parameters  $\boldsymbol{\theta} \in \Theta$  in a  $d$ -dimensional Euclidean space, we wish to determine which  
 1546 hypothesis (possibly with parameter values) is most supported by a sample  $\mathcal{D} = (\omega_1, \dots, \omega_n)^\top$  of  
 1547 materialized outcomes of  $\Omega$ . The Bayes factor  $B_{1,0}$  for  $\mathcal{H}_1$  against  $\mathcal{H}_0$  (*Jeffreys*, 1939; *Kass and Raftery*,  
 1548 1995)

$$1549 \quad B_{1,0} = \frac{p(\mathcal{D}|\mathcal{H}_1)}{p(\mathcal{D}|\mathcal{H}_0)}, \quad (110)$$

1550 summarizes the evidence provided by the data  $\mathcal{D}$  in favor of hypothesis  $\mathcal{H}_1$ , represented by a watershed  
 1551 model, as opposed to the null hypothesis,  $\mathcal{H}_0$ . *Good* (1952) established a simple relationship between the  
 1552 logarithmic score and the logarithmic value of the Bayes factor

$$1553 \quad \log_b(B_{1,0}) = \log_b\left(\frac{p(\mathcal{D}|\mathcal{H}_1)}{p(\mathcal{D}|\mathcal{H}_0)}\right) = S_{\text{LS}}(\mathcal{H}_1, \mathcal{D}) - S_{\text{LS}}(\mathcal{H}_0, \mathcal{D}), \quad (111)$$

1554 where  $\log_b(B_{1,0})$  is also referred as the *weight of evidence*. We can use this identity to compute the values  
 1555 of the Bayes factors for the competing distribution forecasts of Fig. 4a (see Table 14). The listed Bayes  
 1556 factors confirm that  $P_3$  the best predictive density among the distribution forecasts unless, of course, we  
 1557 set  $P = Q$ . The evidence for  $P_3$  is strong as a result of the large sample of ten-thousand observations.

Table 14: Mean values of the Bayes factors  $B_{i,j}$  of the distribution forecasts  $P_1, \dots, P_5$  displayed in Fig. 4a from application of Equation (111) to the sum of the logarithmic scores  $S_{\text{LS}}(P, \omega)$  in units of bits ( $\mathfrak{b} = e$ ) of the  $n = 10^4$  outcomes  $\omega_1, \dots, \omega_n$ . The last row and column correspond to a perfect distribution forecast,  $P = Q$ . *Kass and Raftery (1995)* (P. 777) categorize the Bayes factors and present descriptive statements on the strength of the evidence;  $B_{i,j} = 0$  favors in strongest possible terms the null hypothesis,  $P_j$ , and vice-versa  $B_{i,j} > 75$  means that there is decisive evidence against  $P_j$ .

Forecast	$S_{\text{LS}}(P, Q)$	$P_1$ red	$P_2$ blue	$P_3$ green	$P_4$ yellow	$P_5$ purple	$P = Q$ gray
$P_1$	-6.520		0.000	0.000	0.000	0.000	0.000
$P_2$	-2.422		> 75		0.000	0.000	> 75
$P_3$	-2.190		> 75	> 75		> 75	> 75
$P_4$	-2.260		> 75	> 75	0.000		> 75
$P_5$	-2.479		> 75	0.000	0.000	0.000	
$P = Q$	-2.132		> 75	> 75	> 75	> 75	> 75

1558 For this same reason, we only need to look at the time-averaged values of the logarithmic score  $\bar{S}_{\text{LS}}(P, \omega)$   
1559 in Table 9 to realize that there is overwhelming evidence for the lognormal density.

1560 When the data come in a particular sequence, such as time order, we may develop a more intuitive  
1561 understanding of the above identity if we look at the predictive density of  $\omega_t$  given past observations

1562  $\mathcal{D}^{t-1} = (\omega_1, \dots, \omega_{t-1})^\top$

1563 
$$p(\omega_t | \mathcal{D}^{t-1}, \mathcal{H}) = \int_{\Theta} p(\omega_t | \Theta, \mathcal{H}) p(\Theta | \mathcal{D}^{t-1}, \mathcal{H}) d\Theta. \quad (112)$$

1564 where  $p(\omega_t | \Theta, \mathcal{H})$  is the predictive density of  $\omega_t$  given  $\Theta \in \Theta$ ,  $p(\Theta | \mathcal{D}^{t-1}, \mathcal{H})$  signifies the posterior  
1565 distribution of the parameters and  $t = (1, \dots, n)$ . The logarithmic score for  $\omega_t$  is now equal to

1566 
$$S_{\text{LS}}(\mathcal{H}, \omega_t) = \log_b(p(\omega_t | \mathcal{D}^{t-1}, \mathcal{H})), \quad (113)$$

1567 and we can write

1568 
$$S_{\text{LS}}(\mathcal{H}, \mathcal{D}) = \sum_{t=1}^n \log_b(p(\omega_t | \mathcal{D}^{t-1}, \mathcal{H})), \quad (114)$$

1569 for the total score in Equation (111). The above scoring rule is asymptotically equivalent to the Bayes  
1570 information criterion (*Dawid, 1984*). Thus, the logarithmic score may not only be used to evaluate,  
1571 contrast and rank different distribution forecasts but has broader application to model selection.

#### 1572 8.0.4 Synthesis with model diagnostics

1573 There is an urgent need for *proper* scoring rules of hydrologic signatures in support of diagnostic model  
1574 evaluation. We proposed steps in this direction for the flow duration curve with the CREPS divergence

score but are in need of a much larger family of *strictly proper* signature scoring rules. The relevant number and form of diagnostic signatures is catchment and theory dependent. Let us restrict attention to numeral descriptors of the stream hydrograph, examples of which are the baseflow index, runoff ratio, rising limb density (*Shamir et al.*, 2005) and flashiness index (*Baker et al.*, 2004). Time-averaged values of these indices are frequently used for hydrologic model evaluation, but this projection of the hydrograph to a handful of points implies a significant loss of information. But this loss does not have to be as colossal if we work instead with distribution functions of the hydrologic signatures. We can implement the moving-block bootstrap of *Kunsch* (1989) and obtain frequency distributions of the hydrologic signatures by shifting a window of constant width, say 365 days, by one or more days through the streamflow record and/or hyetograph. The choice of increment exerts control on the smoothness of the signatures' distribution functions. A distribution-based approach to model diagnostics arguably is much more robust and complete in that it, (i) compares measured and simulated signature distributions and not just their mean values, (ii) acknowledges temporal variability of the signatures, and (iii) implicitly accounts for signature uncertainty. The CRPS divergence function in Equation (67) may be used to quantify the statistical distance between measured and simulated signature distributions. Confidence intervals can be derived from the Dvoretzky–Kiefer–Wolfowitz–Massart inequality (*Dvoretzky et al.*, 1956) and its extension by *Naaman* (2021) to multivariate distributions.

The methodology described above evaluates each signature separately, thus, is invariant to the multivariate relationships between the signatures. Suppose  $F_Q$  and  $F_P$  are  $r$ -variate CDFs of the *true* and *simulated* signatures and  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_r)^\top$  is a signature-valued vector drawn at random from  $Q$ . We can generalize distribution-based model diagnostics to the joint distribution of the signatures using the multivariate form of the CRPS

$$S_{\text{MCRPS}}(P, \boldsymbol{\omega}) = - \int_{\Omega} (F_P(\mathbf{u}) - \mathbb{1}\{\boldsymbol{\omega} \leq \mathbf{u}\})^2 d\mathbf{u}, \quad (115)$$

and associated divergence function

$$d_{\text{MCRPS}}(P, Q) = \int_{\Omega} (F_P(\mathbf{u}) - F_Q(\mathbf{u}))^2 d\mathbf{u}, \quad (116)$$

where  $\mathbf{u} \in \Omega \subseteq \mathbb{R}^r$ . By expanding the integrand of Equation (115) we yield a term  $\int_{\Omega} \mathbb{1}\{\boldsymbol{\omega} \leq \mathbf{u}\}^2 d\mathbf{u}$ , which depends only on  $\boldsymbol{\omega}$  and not  $F_P$ . Thus,

$$S_{\text{MCRPS}}^*(P, \boldsymbol{\omega}) = - \int_{\Omega} F_P^2(\mathbf{u}) d\mathbf{u} + 2 \int_{\Omega} F_P(\mathbf{u}) \mathbb{1}\{\boldsymbol{\omega} \leq \mathbf{u}\} d\mathbf{u}, \quad (117)$$

1603 is an affine transformation of the MCRPS scoring rule (see e.g. *Meng et al.* 2022). If the distribution  
 1604  $P$  is made up of  $m$  samples  $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_m\} \in \mathbb{R}^{r \times m}$  of signature values,  $\mathbf{y}_i = (y_{i1}, \dots, y_{ir})^\top \in \mathbb{R}^{r \times 1}$ ,  
 1605 then the multivariate eCDF at point  $\mathbf{u} = (u_1, \dots, u_r)^\top$  is equal to

$$1606 F_P(\mathbf{u}) = F_P(\mathbf{u}|\mathbf{Y}) \triangleq \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{y_{i1} \leq u_1, \dots, y_{ir} \leq u_r\}. \quad (118)$$

1607 The above formulation can be used in Equation (116) for the signatures of the simulated and measured  
 1608 discharge records. *Langrené and Warin* (2021) presents two different algorithms for accurate and CPU-  
 1609 efficient estimation of the multivariate eCDF requiring only  $\mathcal{O}(m)$  operations if samples are sorted. Note  
 1610 that the MCRPS divergence score can be used as a loss function for watershed model parameter estimation  
 1611 so as to maximize the agreement between the measured and simulated joint signature distributions.

1612 The multivariate form of the CRPS divergence can also be used for individual signatures. For example,  
 1613 *Brutsaert and Nieber* (1977), hereafter referred to as BN77, made simplifying assumptions about the  
 1614 catchment water balance in a recession period to arrive at the following relationship between the time  
 1615 rate of change in discharge  $dy/dt$  (mm/d<sup>2</sup>) and discharge  $y$  (mm/d)

$$1616 \frac{dy}{dt} = -ay^b, \quad (119)$$

1617 where  $a$  (d<sup>-1/b</sup>) and  $b$  (-) are unknown recession constants that depend on watershed characteristics. BN77  
 1618 suggest using a log<sub>b</sub> – log<sub>b</sub> graph of  $-dy/dt$  versus  $y$  to estimate  $a$  and  $b$ . This graphical interpretation  
 1619 of recession hydrographs is not without practical problems and has been subject to active debate in the  
 1620 hydrologic literature (*Rupp and Selker*, 2006; *Kirchner*, 2009; *Thomas et al.*, 2013; *Roques et al.*, 2017;  
 1621 *Tashie et al.*, 2020). But, we do not need to know the values of  $a$  and  $b$  for a meaningful model evaluation  
 1622 and/or calibration (e.g. *Jepsen et al.* 2016). We can devise a much stronger test of model performance  
 1623 by comparing directly measured  $Q$  and simulated  $P$  distributions of the  $\log_{10}(y)$  and  $\log_{10}(-dy/dt)$   
 1624 relationship. The divergence of  $P$  and  $Q$  can be measured using the bivariate form of  $d_{\text{MCRPS}}(P, Q)$   
 1625 in Equation (116). The mass-conservative numerical solver of the watershed models will help remedy  
 1626 artefacts of BN77 analysis as a result of a constant time step between successive discharge observations.

1627 To illustrate our idea, we revisit the BMA ensemble and compute the divergence of measured and  
 1628 simulated  $\log_{10}(y)$  and  $\log_{10}(-dy/dt)$  point clouds of the  $K = 8$  conceptual watershed models (see  
 1629 Figure 18). BN77 scatterplot analysis demonstrates that the recession curves of (a) NAM, (b) GR4J, (c)  
 1630 HYMOD and (d) TOPMO are inconsistent with the measured  $y$  and  $-dy/dt$  point cloud. This can be

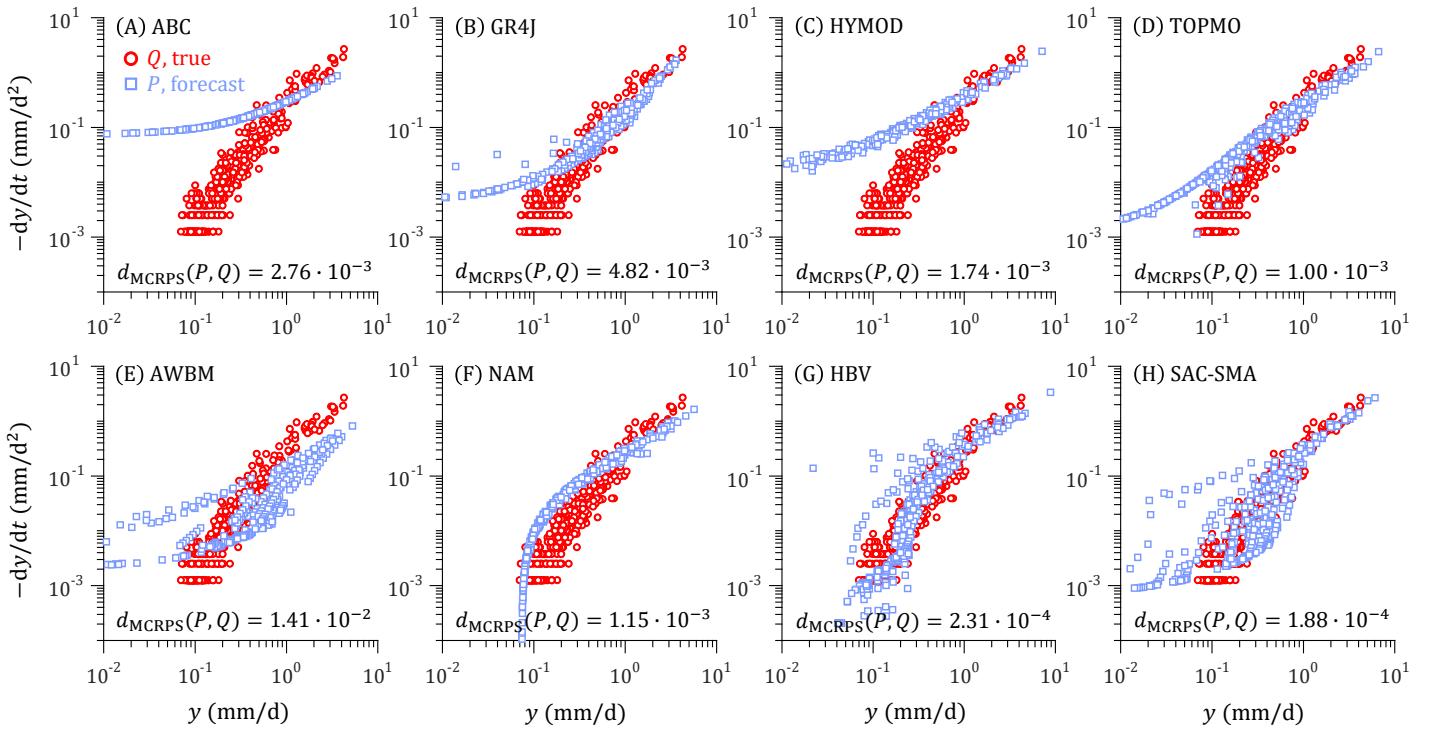


Figure 18: Scatter plots (blue squares) of the  $\log_{10}(y) - \log_{10}(-dy/dt)$  relationship for the different models of the BMA discharge ensemble. The red dots correspond to the measured discharge record.

1631 attributed in part to linear bias-correction, but is reason enough to disqualify these four models from  
 1632 the BMA ensemble. The recession curves of the HBV and SAC-SMA models display the closest match  
 1633 to the measured  $y$  and  $-dy/dt$  point cloud. The MCRPS score divergence,  $d_{\text{MCRPS}}(P, Q)$ , confirms our  
 1634 inherently subjective visual assessment of the watershed models and attains minimum values for the  
 1635 SAC-SMA model and HBV as runner-up.

1636 We can reformulate signatures using existing scoring rules or develop specific scoring rules such as the  
 1637 CREPS divergence for the flow duration curve and MCRPS and CRPS for recession and flood frequency  
 1638 analysis. In general, the extraction of theory-based signatures from measurements of catchment behavior  
 1639 should go hand in hand with the development of adequate scoring rules.

#### 1640 8.0.5 Standardization of model evaluation metrics

1641 The last decades have witnessed an unbridled growth in the number of performance measures used  
 1642 for hydrologic model evaluation. This proliferation is partly a result of the lack of conforming theory  
 1643 and principles for scoring functions. This state of affairs leads to misguided inferences. Elicitability  
 1644 and propriety offer two useful working paradigms for the development, application and evaluation of

1645 hydrologic scoring functions and scoring rules. This will set much higher standards for metric development  
1646 and reduces a model's susceptibility to misinformation and unfinished learning. Note that there is ample  
1647 opportunity to extend hydrologic scoring functions and rules to high-dimensional data from ground-based  
1648 sensor networks and satellites and the spectral domain.

1649 **8.0.6 Practical implementation**

1650 Most hydrologist will be familiar with goodness-of-fit metrics such as the coefficient of determination,  
1651 RMSE, NSE and KGE that quantify the agreement (or distance) between point estimates of measured  
1652 and simulated quantities. Scoring rules, on the contrary, evaluate the accuracy of a distribution forecast  
1653 but may appear more difficult to compute. Therefore, we conclude this paper with a few words about  
1654 their practical use. The case studies of this paper should explicate what is needed for computation of  
1655 the different scoring rules. These requirements are met with a simple GLUE ensemble of simulated  
1656 watershed behavior but are also satisfied with a more sophisticated Bayesian description of predictive  
1657 densities obtained from posterior realizations of the DREAM algorithm (e.g. Vrugt *et al.* 2022). The  
1658 `ScoringRules` toolbox in MATLAB automatically computes the different scoring rules from simulation  
1659 distributions specified in the form of a  $m$ -member ensemble. This toolbox will be released through  
1660 Zenodo upon acceptance.

1661 **9 Conclusions**

1662 This paper was inspired in part by sociocultural developments related to equity, diversity and inclusion  
1663 and advocates the use of simulation distributions for hydrologic model evaluation and diagnostics. A  
1664 simulation distribution coalesces the diversity of model responses across the (prior/posterior) parameter  
1665 and/or input space and contains information about model functioning, behavior, robustness, sensitivity  
1666 and uncertainty that is not available in single-valued model output. But simulation distributions demand  
1667 a fundamentally different approach to model evaluation and diagnostics.

1668 In this paper, we discussed past developments that led to the current state-of-the-art of simulation  
1669 distribution evaluation in hydrology and explored the use of *strictly proper* scoring rules for a more  
1670 robust and complete method of hydrologic model evaluation and diagnostics. *Strictly proper* scoring

rules condense a distribution forecast to a single reward value for the materialized outcome(s) and have a strong underpinning in statistical, decision and information theory. We reviewed scoring rules for dichotomous and categorical events, quantiles (intervals) and density forecasts, discussed the importance of scoring rule propriety and sufficiency, and addressed diagnostic aspects of *strictly proper* scoring rules by means of their analytic decomposition into an entropy (uncertainty), sharpness and reliability term. The power and usefulness of probabilistic (distribution) model evaluation was illustrated by application to a few illustrative examples, 24-h forecasts of daily rainfall, the survival function of discharge and discharge distributions derived from conceptual watershed models using Bayesian model averaging and the GLUE methodology. These studies and the ensuing discussion on future research work, demonstrated the use of scoring rules for diagnostic model evaluation, Bayesian model selection and extreme event prediction. Specifically, we pointed out the relationship between the logarithmic score and Bayes factors, introduced *strictly proper* divergence scores for the flow duration and recession curves and presented an analytic CRPS expression for flood frequency analysis with the Pearson type III distribution.

Our most important results are as follows.

1. Ubiquitous scoring functions such as the mean squared error, coefficient of determination and KG efficiency which quantify model performance using single-valued model output are insensitive to the underlying distribution of this output. The likelihood, on the contrary, is also computed from point-valued model output but preserves information about the underlying distribution of simulated quantities.
2. Distribution forecasts express diversity in the form of a probability distribution over quantities of interests and contain information about model behavior, robustness, sensitivity and uncertainty that is not available in single-valued model output. As we have witnessed in the BMA case study, this can lead to diametrically opposed conclusions about model adequacy.
3. Albeit diagnostically appealing but hydrologic forecast verification metrics such as the coefficient of variation, reliability, width and coverage measure only a few aspects of the statistical consistency between distribution forecasts and verifying observations. This incomplete evaluation does not guarantee honest forecasts.
4. Bregman divergences are of paramount importance in distribution evaluation as they ensure *strict*

1699 *propriety* of the scoring rules and their associated divergence functions. *Strict propriety* implies  
1700 that the score divergence  $d(P, Q)$  of the probabilistic forecast  $P$  and true but unknown distribution  
1701  $Q$  is strictly positive and zero only when  $P = Q$ .

- 1702 5. *Strictly proper* scoring rules such as the quadratic, logarithmic, spherical, continuous ranked  
1703 probability and energy scores facilitate a more robust and complete method of hydrologic model  
1704 evaluation. The larger their values at any time  $t$ , the closer is the simulated distribution to the  
1705 unknown distribution  $Q$  of  $\omega$ .
- 1706 6. The logarithmic score has negative Shannon entropy as its generalized entropy function and the  
1707 relative entropy (reverse KL divergence) as its score divergence.
- 1708 7. The continuous ranked probability score is a generalization of the absolute error (residual) to a  
1709 distribution forecast.
- 1710 8. Propriety and elicability offer two useful working paradigms for the development and application of  
1711 hydrologic scoring functions and scoring rules. Within the context of model evaluation, information-  
1712 theoretic principles set much higher and universal standards for metric development thereby  
1713 promoting metric standardization, reproducibility, stability and comparative analysis across models  
1714 and data sets. Within the context of model training, information-theoretic principles reduce a  
1715 model's susceptibility to misinformation and unfinished learning.
- 1716 9. Sharpness is an *improper* scoring rule as it is a property of the predictive distribution only.
- 1717 10. *Strictly proper* scoring rules can be used to determine the optimum size of a GLUE ensemble.
- 1718 11. The use of time-averaged functionals of the hydrograph such as the baseflow index, runoff ratio,  
1719 rising limb density and flashiness index imply a large loss of information about the temporal  
1720 behavior of watershed functioning.
- 1721 12. Frequency distributions of numeral streamflow descriptors derived from a moving-block bootstrap  
1722 method enable a more robust and complete model evaluation and diagnostics by means of *strictly*  
1723 *proper* scoring rules. The signature frequency distributions may be treated separately or as a joint  
1724 multivariate distribution.

- 1725 13. The CREPS divergence measures in a single reward value the agreement of measured  $R_Q$  and  
1726 simulated  $R_P$  flow duration curves. This function is strictly positive and zero only if  $R_P = R_Q$ ,  
1727 expressed in physical units of discharge, sensitive to all moments of the FDC and equal to the  
1728 distance (absolute error) between two point measures of the FDC.
- 1729 14. The bivariate form of the continuous ranked probability score offers a *strictly proper* scoring rule  
1730 for BN77 recession analysis of the joint distribution of the time rate of change in discharge and  
1731 discharge itself. The comparison of observed and simulated point clouds simplifies model diagnostics  
1732 and selection.
- 1733 15. Differential weighting provides an opportunity to hone in on specific aspects of the simulation  
1734 distribution. These so-called local and/or censored scoring functions do not sacrifice local propriety  
1735 and are specifically suited to evaluating a models' ability in simulating accurately the frequency of  
1736 extreme events.
- 1737 16. There is a pressing need for *proper* and *strictly proper* scoring rules of hydrologic signatures. The  
1738 CREPS divergence of the flow duration curve is an important step in this direction but model  
1739 diagnostics would benefit from decision-theoretically principled hydrograph functionals.
- 1740 Last but not least, but currently, there is no delineation or field of science in hydrology and water  
1741 resources that expects and appreciates mathematical and statistical rigor. By analogy with the fields of  
1742 biometrics, econometrics and psychometrics, which grew out of the application of theory, mathematics,  
1743 and statistical inference to biology, economy and psychology, respectively, we believe that this work would  
1744 fit naturally under *hydrometrics*. Though, this requires a redefinition of hydrometry by the International  
1745 Organization for Standardization “...the science of monitoring water in natural water resources” to  
1746 include the application of theory, mathematics and statistical inference to study, model and predict the  
1747 water distribution and movement. This would develop the new title of hydrometrician.

## 1748 Data and Software Availability

- 1749 The different case studies presented in this paper are part of the MATLAB toolbox `Scoring_Rules`.  
1750 This toolbox will be made available for download from Zenodo upon acceptance. Data, models and

<sub>1751</sub> algorithms are organized in separate folders. The MODELAVG toolbox of (*Vrugt and Beven, 2018*) can  
<sub>1752</sub> be obtained from GitHub at <https://github.com/jaspervrugt/MODELAVG>. Hydrologic data from  
<sub>1753</sub> the CAMELS data set is described in *Newman et al. (2015)* and can be downloaded from <https://dx.doi.org/10.5065/D6MW2F4D>.  
<sub>1754</sub>

1755    **References**

- 1756    1. Alexander, C., M. Coulon, Y. Han, and X. Meng (2022), Evaluating the discrimination ability of  
1757    proper multi-variate scoring rules, *Annals of Operations Research*, doi:10.1007/s10479-022-04611-9.
- 1758    2. Ammann, L., F. Fenicia, and P. Reichert (2019), A likelihood framework for deterministic  
1759    hydrological models and the importance of non-stationary autocorrelation, *Hydrology and Earth  
1760    System Sciences*, 23(4), 2147–2172, doi:10.5194/hess-23-2147-2019.
- 1761    3. Baker, D. B., R. P. Richards, T. T. Loftus, and J. W. Kramer (2004), A new flashiness index:  
1762    Characteristics and applications to midwestern rivers and streams, *JAWRA Journal of the  
1763    American Water Resources Association*, 40(2), 503–522, doi:10.1111/j.1752-1688.2004.tb01046.x.
- 1764    4. Baringhaus, L., and C. Franz (2004), On a new multivariate two-sample test, *Journal of Multi-  
1765    variate Analysis*, 88, 190–206, doi:10.1016/S0047-259X(03)00079-4.
- 1766    5. Bates, B. C., and E. P. Campbell (2001), A Markov chain Monte Carlo scheme for parameter  
1767    estimation and inference in conceptual rainfall-runoff modeling, *Water Resources Research*, 37(4),  
1768    937–947, doi:10.1029/2000WR900363.
- 1769    6. Bayes, T. (1763), An essay toward solving a problem in the doctrine of chances. By the late Rev.  
1770    Mr. Bayes, F. R. S. communicated by Mr. Price, in a letter to John Canton, A. M. F. R. S.,  
1771    *Philosophical Transactions of the Royal Society of London*, 53, 370–418, doi:10.1098/rstl.1763.0053.
- 1772    7. Bella, A., C. Ferri, J. Hernández-Orallo, Ramírez-Quintana, and J. María (2013), On the  
1773    effect of calibration in classifier combination, *Applied Intelligence*, 38, 566–585, doi:10.1007/  
1774    s10489-012-0388-2.
- 1775    8. Bernardo, J. M. (1979), Expected information as expected utility, *The Annals of Statistics*, 7,  
1776    686–690, doi:10.1214/aos/1176344689.
- 1777    9. Beven, K. (2006), A manifesto for the equifinality thesis, *Journal of Hydrology*, 320(1), 18–36,  
1778    doi:10.1016/j.jhydrol.2005.07.007.
- 1779    10. Beven, K., and A. Binley (1992), The future of distributed models: Model calibration and  
1780    uncertainty prediction, *Hydrological Processes*, 6(3), 279–298, doi:10.1002/hyp.3360060305.

- 1781 11. Beven, K., and J. Freer (2001), Equifinality, data assimilation, and uncertainty estimation in  
1782 mechanistic modelling of complex environmental systems using the GLUE methodology, *Journal*  
1783 of Hydrology, 249(1), 11–29, doi:10.1016/S0022-1694(01)00421-8.
- 1784 12. Bouhlel, N., and A. Dziri (2019), Kullback–leibler divergence between multivariate generalized  
1785 gaussian distributions, *IEEE Signal Processing Letters*, 26(7), 1021–1025, doi:10.1109/LSP.2019.  
1786 2915000.
- 1787 13. Boyle, D. P. (2001), *Multicriteria calibration of hydrological models (PhD thesis)*, Department of  
1788 Hydrology and Water Resources, University of Arizona, Tucson, AZ.
- 1789 14. Boyle, D. P., H. V. Gupta, and S. Sorooshian (2000), Toward improved calibration of hydrologic  
1790 models: Combining the strengths of manual and automatic methods, *Water Resources Research*,  
1791 36(12), 3663–3674, doi:10.1029/2000WR900207.
- 1792 15. Bracher, J., E. L. Ray, T. Gneiting, and N. G. Reich (2021), Evaluating epidemic forecasts in an  
1793 interval format, *PLOS Computational Biology*, 17(2), 1–15, doi:10.1371/journal.pcbi.1008618.
- 1794 16. Bregman, L. (1967), The relaxation method of finding the common point of convex sets and its  
1795 application to the solution of problems in convex programming, *USSR Computational Mathematics*  
1796 and *Mathematical Physics*, 7(3), 200–217, doi:10.1016/0041-5553(67)90040-7.
- 1797 17. Brier, G. W. (1950), Verification of forecasts expressed in terms of probability, *Monthly Weather  
1798 Review*, 78(1), 1–3.
- 1799 18. Bröcker, J. (2009), Reliability, sufficiency, and the decomposition of proper scores, *Quarterly  
1800 Journal of the Royal Meteorological Society*, 135(643), 1512–1519, doi:10.1002/qj.456.
- 1801 19. Brunetti, C., N. Linde, and J. A. Vrugt (2017), Bayesian model selection in hydrogeophysics:  
1802 Application to conceptual subsurface models of the south oyster bacterial transport site, virginia,  
1803 usa, *Advances in Water Resources*, 102, 127–141, doi:10.1016/j.advwatres.2017.02.006.
- 1804 20. Brutsaert, W., and J. L. Nieber (1977), Regionalized drought flow hydrographs from a mature  
1805 glaciated plateau, *Water Resources Research*, 13(3), 637–643, doi:10.1029/WR013i003p00637.

- 1806 21. Buja, A., W. Stuetzle, and Y. Shen (2005), Loss Functions for Binary Class Probability Estimation  
1807 and Classification: Structure and Applications, *Tech. rep.*, Statistics Department, The Wharton  
1808 School, University of Pennsylvania, Philadelphia, PA 19104-6302.
- 1809 22. Burnash, R. J. C., R. L. Ferral, and R. A. McGuire (1973), A generalized streamflow simulation  
1810 system: Conceptual modeling for digital computers, *Tech. rep.*, Joint Federal-State River Forecast  
1811 Center: US Department of Commerce, National Weather Service and CA Department of Water  
1812 Resources, Sacramento, CA.
- 1813 23. Burnham, K. P., and D. R. Anderson (2002), *Model Selection and Multimodel Inference: A  
1814 Practical Information-Theoretic Approach*, 2 ed., 488 pp., Springer, New York, doi:10.1007/b97636.
- 1815 24. Casella, G., and R. L. Berger (2002), *Statistical Inference*, Duxbury Advanced Series, 2 ed., 660  
1816 pp., Duxbury, Pacific Grove, CA.
- 1817 25. Cervera, J. L., and J. Muñoz (1996), Proper scoring rules for fractiles, in *Bayesian Statistics 5*,  
1818 edited by J. M. Bernardo, J. O. Berger, A. P. Dawid, and A. F. M. Smith, pp. 513–519, Oxford  
1819 University Press.
- 1820 26. Chang, H., Y. Yao, A. Koschan, B. Abidi, and M. Abidi (2009), Improving face recognition via  
1821 narrowband spectral range selection using jeffrey divergence, *IEEE Transactions on Information  
1822 Forensics and Security*, 4(1), 111–122, doi:10.1109/TIFS.2008.2012211.
- 1823 27. Chen, L., Z. Shen, X. Yang, Q. Liao, and S. L. Yu (2014), An interval-deviation approach for  
1824 hydrology and water quality model evaluation within an uncertainty framework, *Journal of  
1825 Hydrology*, 509, 207–214, doi:10.1016/j.jhydrol.2013.11.043.
- 1826 28. Christensen, H. M. (2015), Decomposition of a new proper score for verification of ensemble  
1827 forecasts, *Monthly Weather Review*, 143(5), 1517–1532, doi:10.1175/MWR-D-14-00150.1.
- 1828 29. Christoffersen, P. F. (1998), Evaluating interval forecasts, *International Economic Review*, 39(4),  
1829 841–862.
- 1830 30. Clark, M. P., A. G. Slater, D. E. Rupp, R. A. Woods, J. A. Vrugt, H. V. Gupta, T. Wagener, and  
1831 L. E. Hay (2008), Framework for understanding structural errors (FUSE): A modular framework

- 1832 to diagnose differences between hydrological models, *Water Resources Research*, 44(12), doi:  
1833 10.1029/2007WR006735.
- 1834 31. Cover, T. M., and J. A. Thomas (2006), *Elements of Information Theory*, Telecommunications  
1835 and Signal Processing, 2 ed., 784 pp., John Wiley & Sons, Inc., Hoboken, NJ.
- 1836 32. Cox, D. R., and P. A. W. Lewis (1966), *Statistical Analysis of Series of Events*, Methuen's  
1837 Monographs on Applied Probability and Statistics (MMAPS), 285 pp., Methuen, London.
- 1838 33. Csiszar, I. (1975), *I-Divergence Geometry of Probability Distributions and Minimization Problems*,  
1839 *The Annals of Probability*, 3(1), 146–158, doi:10.1214/aop/1176996454.
- 1840 34. Dawid, A. P. (1984), Present position and potential developments: Some personal views: Statistical  
1841 theory: The prequential approach, *Journal of the Royal Statistical Society. Series A (General)*,  
1842 147(2), 278–292.
- 1843 35. Dawid, A. P. (1998), Coherent measures of discrepancy, uncertainty and dependence, with  
1844 applications to Bayesian predictive experimental design, *Tech. rep.*, Department of Statistical  
1845 Science, University College London, London, UK.
- 1846 36. Dawid, A. P., and M. Musio (2014), Theory and applications of proper scoring rules, *METRON*,  
1847 72, 169–183, doi:10.1007/s40300-014-0039-y.
- 1848 37. Dawid, P., and P. Sebastiani (1999), Coherent dispersion criteria for optimal experimental design,  
1849 *The Annals of Statistics*, 27(1), 65–81, doi:10.1214/aos/1018031101.
- 1850 38. de Finetti, B. (2017), *Theory of Probability: A Critical Introductory Treatment*, Wiley Series in  
1851 Probability and Statistics, John Wiley & Sons Ltd, United Kingdom.
- 1852 39. de Punder, R., C. Diks, R. Laeven, and D. van Dijk (2022), A general procedure for localising  
1853 strictly proper scoring rules, *Tech. rep.*, Erasmus University Rotterdam, Tinbergen Institute,  
1854 Rotterdam, The Netherlands.
- 1855 40. DeGroot, M. H., and S. E. Fienberg (1983), The comparison and evaluation of forecasters, *Journal*  
1856 *of the Royal Statistical Society. Series D (The Statistician)*, 32(1/2), 12–22.

- 1857 41. Diebold, F. X., T. A. Gunther, and A. S. Tay (1998), Evaluating density forecasts with applications  
1858 to financial risk management, *International Economic Review*, 39(4), 863–883.
- 1859 42. Diks, C., V. Panchenko, and D. van Dijk (2011), Likelihood-based scoring rules for comparing  
1860 density forecasts in tails, *Journal of Econometrics*, 163(2), 215–230, doi:10.1016/j.jeconom.2011.  
1861 04.001.
- 1862 43. Dimitriadis, T., T. Gneiting, and A. I. Jordan (2021), Stable reliability diagrams for probabilistic  
1863 classifiers, *Proceedings of the National Academy of Sciences*, 118(8), e2016191,118, doi:10.1073/  
1864 pnas.2016191118.
- 1865 44. Duffie, D., and J. Pan (1997), An overview of value at risk, *The Journal of Derivatives*, 4(3),  
1866 7–49, doi:10.3905/jod.1997.407971.
- 1867 45. Dunsmore, I. R. (1968), A Bayesian approach to calibration, *Journal of the Royal Statistical  
1868 Society. Series B (Methodological)*, 30(2), 396–405.
- 1869 46. Dvoretzky, A., J. Kiefer, and J. Wolfowitz (1956), Asymptotic Minimax Character of the Sample  
1870 Distribution Function and of the Classical Multinomial Estimator, *The Annals of Mathematical  
1871 Statistics*, 27(3), 642–669, doi:10.1214/aoms/1177728174.
- 1872 47. Evin, G., D. Kavetski, M. Thyre, and G. Kuczera (2013), Pitfalls and improvements in the  
1873 joint inference of heteroscedasticity and autocorrelation in hydrological model calibration, *Water  
1874 Resources Research*, 49(7), 4518–4524, doi:10.1002/wrcr.20284.
- 1875 48. Evin, G., M. Thyre, D. Kavetski, D. McInerney, and G. Kuczera (2014), Comparison of joint  
1876 versus postprocessor approaches for hydrological uncertainty estimation accounting for error  
1877 autocorrelation and heteroscedasticity, *Water Resources Research*, 50(3), 2350–2375, doi:10.1002/  
1878 2013WR014185.
- 1879 49. Fenicia, F., D. Kavetski, and H. H. G. Savenije (2011), Elements of a flexible approach for  
1880 conceptual hydrological modeling: 1. motivation and theoretical development, *Water Resources  
1881 Research*, 47(11), doi:10.1029/2010WR010174.

- 1882 50. Ferson, S., W. L. Oberkampf, and L. Ginzburg (2008), Model validation and predictive capability  
1883 for the thermal challenge problem, *Computer Methods in Applied Mechanics and Engineering*,  
1884 197(29), 2408–2430, doi:10.1016/j.cma.2007.07.030, validation Challenge Workshop.
- 1885 51. Freer, J., K. Beven, and B. Ambroise (1996), Bayesian estimation of uncertainty in runoff prediction  
1886 and the value of data: An application of the GLUE approach, *Water Resources Research*, 32(7),  
1887 2161–2173, doi:10.1029/95WR03723.
- 1888 52. Friederichs, P., and A. Hense (2007), Statistical downscaling of extreme precipitation events using  
1889 censored quantile regression, *Monthly Weather Review*, 135(6), 2365–2378, doi:10.1175/MWR3403.  
1890 1.
- 1891 53. Friederichs, P., and T. L. Thorarinsdottir (2012), Forecast verification for extreme value distribu-  
1892 tions with an application to probabilistic peak wind prediction, *Environmetrics*, 23(7), 579–594,  
1893 doi:10.1002/env.2176.
- 1894 54. Friedman, D. (1983), Effective scoring rules for probabilistic forecasts, *Management Science*, 29(4),  
1895 447–454, doi:10.1287/mnsc.29.4.447.
- 1896 55. Garratt, A., K. Lee, M. H. Pesaran, and Y. Shin (2003), Forecast uncertainties in macroeconomic  
1897 modeling: An application to the u.k. economy, *Journal of the American Statistical Association*,  
1898 98(464), 829–838.
- 1899 56. Garthwaite, P. H., J. B. Kadane, and A. O'Hagan (2005), Statistical methods for eliciting  
1900 probability distributions, *Journal of the American Statistical Association*, 100(470), 680–701,  
1901 doi:10.1198/016214505000000105.
- 1902 57. Gerstner, T., and M. Griebel (1998), Numerical integration using sparse grids, *Numerical Algo-*  
1903 *rithms*, 18, 209–232, doi:10.1023/A:1019129717644.
- 1904 58. Gerstner, T., and M. Griebel (2003), Dimension-adaptive tensor-product quadrature, *Computing*,  
1905 71(1), 65–87, doi:10.1007/s00607-003-0015-5.
- 1906 59. Gini, C. (1909), Concentration and dependency ratios (in italian), *Rivista di Politica Economica*,  
1907 87(8-9), 769–790.

- 1908 60. Girons Lopez, M., L. Crochemore, and I. G. Pechlivanidis (2021), Benchmarking an operational  
1909 hydrological model for providing seasonal forecasts in sweden, *Hydrology and Earth System*  
1910 *Sciences*, 25(3), 1189–1209, doi:10.5194/hess-25-1189-2021.
- 1911 61. Gneiting, T. (2011), Making and evaluating point forecasts, *Journal of the American Statistical*  
1912 *Association*, 106(494), 746–762, doi:10.1198/jasa.2011.r10138.
- 1913 62. Gneiting, T., and A. E. Raftery (2007), Strictly proper scoring rules, prediction, and es-  
1914 timation, *Journal of the American Statistical Association*, 102(477), 359–378, doi:10.1198/  
1915 016214506000001437.
- 1916 63. Gneiting, T., and R. Ranjan (2011), Comparing density forecasts using threshold-and quantile-  
1917 weighted scoring rules, *Journal of Business & Economic Statistics*, 29(3), 411–422.
- 1918 64. Gneiting, T., F. Balabdaoui, and A. E. Raftery (2005), Probabilistic forecasts, calibration, and  
1919 sharpness, *Tech. rep.*, Department of Statistics, University of Washington, Seattle, WA, USA.
- 1920 65. Gneiting, T., F. Balabdaoui, and A. E. Raftery (2007), Probabilistic forecasts, calibration and  
1921 sharpness, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(2),  
1922 243–268, doi:10.1111/j.1467-9868.2007.00587.x.
- 1923 66. Gong, W., H. V. Gupta, D. Yang, K. Sricharan, and A. O. Hero III (2013), Estimating epistemic  
1924 and aleatory uncertainties during hydrologic modeling: An information theoretic approach, *Water*  
1925 *Resources Research*, 49(4), 2253–2273, doi:10.1002/wrcr.20161.
- 1926 67. Good, I. J. (1952), Rational decisions, *Journal of the Royal Statistical Society, Series B (Statistical*  
1927 *Methodology)*, 14(1), 107–114.
- 1928 68. Granger, C. W. J. (2005), Preface: Some thoughts on the future of forecasting, *Oxford Bulletin of*  
1929 *Economics and Statistics*, 67(s1), 707–711, doi:10.1111/j.1468-0084.2005.00138.x.
- 1930 69. Grayson, R., and G. Blöschl (2001), *Spatial Patterns in Catchment Hydrology: Observations and*  
1931 *Modelling*, p. 416, Cambridge University Press.

- 1932 70. Grimit, E. P., T. Gneiting, B. V. J., and N. A. Johnson (2006), The continuous ranked probability  
1933 score for circular variables and its application to mesoscale forecast ensemble verification, *Quarterly*  
1934 *Journal of the Royal Meteorological Society*, 132(621C), 2925–2942, doi:10.1256/qj.05.235.
- 1935 71. Grünwald, P. D., and A. P. Dawid (2004), Game theory, maximum entropy, minimum discrepancy  
1936 and robust Bayesian decision theory, *The Annals of Statistics*, 32(4), 1367 – 1433, doi:10.1214/  
1937 009053604000000553.
- 1938 72. Grushka-Cockayne, Y., K. C. Lichtendahl, V. R. R. Jose, and R. L. Winkler (2017), Quantile  
1939 evaluation, sensitivity to bracketing, and sharing business payoffs, *Operations Research*, 65(3),  
1940 712–728, doi:10.1287/opre.2017.1588.
- 1941 73. Gupta, H. V., S. Sorooshian, and P. O. Yapo (1998), Toward improved calibration of hydrologic  
1942 models: Multiple and noncommensurable measures of information, *Water Resources Research*,  
1943 34(4), 751–763, doi:10.1029/97WR03495.
- 1944 74. Gupta, H. V., T. Wagener, and Y. Liu (2008), Reconciling theory with observations: elements  
1945 of a diagnostic approach to model evaluation, *Hydrological Processes*, 22(18), 3802–3813, doi:  
1946 10.1002/hyp.6989.
- 1947 75. Gupta, H. V., H. Kling, K. K. Yilmaz, and G. F. Martinez (2009), Decomposition of the mean  
1948 squared error and NSE performance criteria: Implications for improving hydrological modelling,  
1949 *Journal of Hydrology*, 377(1), 80–91, doi:10.1016/j.jhydrol.2009.08.003.
- 1950 76. Guttorp, P. (2011), The role of statisticians in international science policy, *Environmetrics*, 22(7),  
1951 817–825, doi:10.1002/env.1109.
- 1952 77. Halton, J. H. (1960), On the efficiency of certain quasi-random sequences of points in evaluating  
1953 multidimensional integrals, *Numerische Mathematik*, 2, 84–90, doi:10.1007/BF01386213.
- 1954 78. Hamill, T. M. (2001), Interpretation of rank histograms for verifying ensemble forecasts, *Monthly*  
1955 *Weather Review*, 129(3), 550–560, doi:10.1175/1520-0493(2001)129<0550:IORHFV>2.0.CO;2.
- 1956 79. Hammersley, J. M., and D. C. Handscomb (1960), *Monte Carlo Methods*, Monographs on  
1957 Statistics and Applied Probability (MSAP), 1 ed., 178 pp., Springer, Dordrecht, doi:10.1007/  
1958 978-94-009-5819-7.

- 1959 80. Hersbach, H. (2000), Decomposition of the continuous ranked probability score for ensemble pre-  
1960 diction systems, *Weather and Forecasting*, 15(5), 559–570, doi:10.1175/1520-0434(2000)015<0559:  
1961 DOTCRP>2.0.CO;2.
- 1962 81. Hosking, J. R. M., and J. R. Wallis (1997), *Regional Frequency Analysis: An Approach Based on*  
1963 *L-moments*, p. 224, Cambridge University Press, doi:10.1017/cbo9780511529443.
- 1964 82. Hughes, G., and C. F. Topp (2015), Probabilistic forecasts: Scoring rules and their decomposition  
1965 and diagrammatic representation via Bregman divergences, *Entropy*, 17(8), 5450–5471, doi:  
1966 10.3390/e17085450.
- 1967 83. Hunt, V., S. Dixon-Fyle, S. Prince, and K. Dolan (2020), Diversity wins: How inclusion matters,  
1968 *Tech. rep.*, McKinsey & Company.
- 1969 84. Jakeman, J. D., and A. Narayan (2018), Generation and application of multivariate polynomial  
1970 quadrature rules, *Computer Methods in Applied Mechanics and Engineering*, 338, 134–161, doi:  
1971 10.1016/j.cma.2018.04.009.
- 1972 85. Jaynes, E. T. (1963), *Information Theory and Statistical Mechanics*, pp. 181–218, W. A. Benjamin,  
1973 Inc.
- 1974 86. Jeffreys, H. (1939), *The theory of probability*, p. 470, Oxford Classic Texts in the Physical Sciences,  
1975 3 ed., Oxford University Press.
- 1976 87. Jeffreys, H. (1946), An invariant form for the prior probability in estimation problems, *Proceedings*  
1977 *of the Royal Society of London. Series A, Mathematical and Physical Sciences*, 186(1007), 453–461.
- 1978 88. Jepsen, S. M., T. C. Harmon, and Y. Shi (2016), Watershed model calibration to the base flow  
1979 recession curve with and without evapotranspiration effects, *Water Resources Research*, 52(4),  
1980 2919–2933, doi:10.1002/2015WR017827.
- 1981 89. Jolliffe, I. T., and D. B. Stephenson (2011), *Forecast Verification: A Practitioner's Guide in*  
1982 *Atmospheric Science*, p. 296, 2 ed., Wiley Blackwell.
- 1983 90. Jordan, A. (2016), Facets of forecast evaluation, Ph.D. thesis, Karlsruhe Institute of Technology,  
1984 doi:10.5445/IR/1000063629.

- 1985 91. Kass, R. E., and A. E. Raftery (1995), Bayes factors, *Journal of the American Statistical  
Association*, 90(430), 773–795, doi:10.1080/01621459.1995.10476572.
- 1986
- 1987 92. Kavetski, D., G. Kuczera, and S. W. Franks (2006a), Bayesian analysis of input uncertainty in  
hydrological modeling: 1. Theory, *Water Resources Research*, 42(3), doi:10.1029/2005WR004368.
- 1988
- 1989 93. Kavetski, D., G. Kuczera, and S. W. Franks (2006b), Bayesian analysis of input uncertainty in hy-  
drological modeling: 2. Application, *Water Resources Research*, 42(3), doi:10.1029/2005WR004376.
- 1990
- 1991 94. Kirchner, J. W. (2009), Catchments as simple dynamical systems: Catchment characterization,  
rainfall-runoff modeling, and doing hydrology backward, *Water Resources Research*, 45(2), doi:  
10.1029/2008WR006912.
- 1992
- 1993
- 1994 95. Knoben, W. J. M., J. E. Freer, K. J. A. Fowler, M. C. Peel, and R. A. Woods (2019), Modular  
assessment of rainfall–runoff models toolbox (marrmot) v1.2: an open-source, extendable frame-  
work providing implementations of 46 conceptual hydrologic models as continuous state-space  
formulations, *Geoscientific Model Development*, 12(6), 2463–2480, doi:10.5194/gmd-12-2463-2019.
- 1995
- 1996
- 1997
- 1998 96. Knorr-Held, L., and E. Rainer (2001), Projections of lung cancer in West Germany: A case study  
in Bayesian prediction, *Biostatistics*, 2(1), 109–129, doi:10.1093/biostatistics/2.1.109.
- 1999
- 2000 97. Kolmogorov, A. N. (1933), Sulla determinazione empirica di una legge di distribuzione, *Giornale  
dell'Istituto Italiano degli Attuari*, 4, 83–91.
- 2001
- 2002 98. Kuczera, G. (1983), Improved parameter inference in catchment models: 1. Evaluating parameter  
uncertainty, *Water Resources Research*, 19(5), 1151–1162, doi:10.1029/WR019i005p01151.
- 2003
- 2004 99. Kuczera, G., and E. Parent (1998), Monte Carlo assessment of parameter uncertainty in conceptual  
catchment models: the Metropolis algorithm, *Journal of Hydrology*, 211(1), 69–85, doi:10.1016/S0022-1694(98)00198-X.
- 2005
- 2006
- 2007 100. Kull, M., and P. Flach (2015), Novel decompositions of proper scoring rules for classification:  
Score adjustment as precursor to calibration, in *Machine Learning and Knowledge Discovery  
in Databases*, edited by A. Appice, P. P. Rodrigues, V. Santos Costa, C. Soares, J. Gama, and  
A. Jorge, pp. 68–85, Springer International Publishing, Cham.
- 2008
- 2009
- 2010

- 2011 101. Kullback, S. (1959), *Information Theory and Statistics*, John Wiley & Sons, New York.
- 2012 102. Kullback, S., and R. A. Leibler (1951), On information and sufficiency, *Annals of Mathematical*  
2013 *Statistics*, 22(1), 79–86, doi:10.1214/aoms/1177729694.
- 2014 103. Kunsch, H. R. (1989), The Jackknife and the Bootstrap for General Stationary Observations, *The*  
2015 *Annals of Statistics*, 17(3), 1217–1241, doi:10.1214/aos/1176347265.
- 2016 104. Lai, T. L., S. T. Gross, and D. B. Shen (2011), Evaluating probability forecasts, *The Annals of*  
2017 *Statistics*, 39(5), 2356–2382.
- 2018 105. Laio, F., and S. Tamea (2007), Verification tools for probabilistic forecasts of continuous  
2019 hydrological variables, *Hydrology and Earth System Sciences*, 11(4), 1267–1277, doi:10.5194/  
2020 hess-11-1267-2007.
- 2021 106. Lamontagne, J. R., C. A. Barber, and R. M. Vogel (2020), Improved estimators of model  
2022 performance efficiency for skewed hydrologic data, *Water Resources Research*, 56(9), 1–25, doi:  
2023 10.1029/2020wr027101.
- 2024 107. Langrené, N., and X. Warin (2021), Fast multivariate empirical cumulative distribution function  
2025 with connection to kernel density estimation, *Computational Statistics & Data Analysis*, 162,  
2026 107,267, doi:10.1016/j.csda.2021.107267.
- 2027 108. Liu, Y., W. Chen, P. Arendt, and H.-Z. Huang (2011), Toward a better understanding of model  
2028 validation metrics, *Journal of Mechanical Design*, 133, 071,005.
- 2029 109. Lu, D., M. Ye, and S. P. Neuman (2011), Dependence of bayesian model selection criteria  
2030 and fisher information matrix on sample size, *Mathematical Geosciences*, 43, 971–993, doi:  
2031 10.1007/s11004-011-9359-0.
- 2032 110. Luke, A., J. A. Vrugt, A. AghaKouchak, R. Matthew, and B. F. Sanders (2017), Predicting  
2033 nonstationary flood frequencies: Evidence supports an updated stationarity thesis in the united  
2034 states, *Water Resources Research*, 53(7), 5469–5494, doi:10.1002/2016WR019676.
- 2035 111. Matheron, G. (1984), *The Selectivity of the Distributions and the "Second Principle of Geostatistics"*,  
2036 p. 421–433, Springer, doi:10.1007/978-94-009-3699-7\_24.

- 2037 112. Matheson, J. E., and R. L. Winkler (1976), Scoring rules for continuous probability distributions,  
2038     *Management Science*, 22(10), 1087–1096, doi:10.1287/mnsc.22.10.1087.
- 2039 113. McCarthy, J. (1956), Measures of the value of information, *Proceedings of the National Academy  
2040     of Sciences*, 42(9), 654–655, doi:10.1073/pnas.42.9.654.
- 2041 114. McDonald, J. B., and B. C. Jensen (1979), An analysis of some properties of alternative measures  
2042     of income inequality based on the gamma distribution function, *Journal of the American Statistical  
2043     Association*, 74(368), 856–860.
- 2044 115. McInerney, D., M. Thyre, D. Kavetski, J. Lerat, and G. Kuczera (2017), Improving probabilistic  
2045     prediction of daily streamflow by identifying pareto optimal approaches for modeling heteroscedastic  
2046     residual errors, *Water Resources Research*, 53(3), 2199–2239, doi:10.1002/2016WR019168.
- 2047 116. McInerney, D., D. Kavetski, M. Thyre, J. Lerat, and G. Kuczera (2019), Benefits of explicit  
2048     treatment of zero flows in probabilistic hydrological modeling of ephemeral catchments, *Water  
2049     Resources Research*, 55(12), 11,035–11,060, doi:10.1029/2018WR024148.
- 2050 117. Meng, X., J. W. Taylor, S. Ben Taieb, and S. Li (2022), Scores for multivariate distributions and  
2051     level sets, *Tech. rep.*, University of Sussex, Falmer, Brighton, UK, BN1 9SN, doi:10.48550/arXiv.  
2052     2002.09578.
- 2053 118. Murphy, A. H. (1970), The ranked probability score and the probability score: A comparison,  
2054     *Monthly Weather Review*, 98(12), 917–924, doi:10.1175/1520-0493(1970)098<0917:TRPSAT>2.3.  
2055     CO;2.
- 2056 119. Murphy, A. H. (1973), A new vector partition of the probability score, *Journal of Applied  
2057     Meteorology and Climatology*, 12(4), 595 – 600, doi:10.1175/1520-0450(1973)012<0595:ANVPOT>  
2058     2.0.CO;2.
- 2059 120. Murphy, A. H. (1993), What is a good forecast? an essay on the nature of goodness in weather  
2060     forecasting, *Weather and Forecasting*, 8(2), 281–293, doi:10.1175/1520-0434(1993)008<0281:  
2061     WIAGFA>2.0.CO;2.

- 2062 121. Murphy, A. H. (1996), General decompositions of mse-based skill scores: Measures of some  
2063 basic aspects of forecast quality, *Monthly Weather Review*, 124(10), 2353–2369, doi:10.1175/  
2064 1520-0493(1996)124<2353:GDOMBS>2.0.CO;2.
- 2065 122. Murphy, A. H., and R. W. Katz (1985), *Probability, Statistics, and Decision Making in the*  
2066 *Atmospheric Sciences*, p. 560, 1 ed., Westview Press.
- 2067 123. Murphy, A. H., and R. L. Winkler (1987), A general framework for forecast verification, *Monthly*  
2068 *Weather Review*, 115(7), 1330–1338, doi:10.1175/1520-0493(1987)115<1330:AGFFFV>2.0.CO;2.
- 2069 124. Naaman, M. (2021), On the tight constant in the multivariate dvoretzky–kiefer–wolfowitz inequality,  
2070 *Statistics & Probability Letters*, 173, 109,088, doi:10.1016/j.spl.2021.109088.
- 2071 125. Nash, J., and J. Sutcliffe (1970), River flow forecasting through conceptual models part I - A  
2072 discussion of principles, *Journal of Hydrology*, 10(3), 282–290, doi:10.1016/0022-1694(70)90255-6.
- 2073 126. Neuman, S. P. (2003), Maximum likelihood bayesian averaging of uncertain model predic-  
2074 tions, *Stochastic Environmental Research and Risk Assessment*, 17, 291–305, doi:10.1007/  
2075 s00477-003-0151-7.
- 2076 127. Newman, A. J., et al. (2015), Development of a large-sample watershed-scale hydrometeorological  
2077 data set for the contiguous USA: data set characteristics and assessment of regional variability  
2078 in hydrologic model performance, *Hydrology and Earth System Sciences*, 19(1), 209–223, doi:  
2079 10.5194/hess-19-209-2015.
- 2080 128. Niederreiter, H. (1992), *Random Number Generation and Quasi-Monte Carlo Methods*, CBMS-NSF  
2081 Regional Conference Series in Applied Mathematics, 237 pp., Society for Industrial and Applied  
2082 Mathematics, doi:10.1137/1.9781611970081.
- 2083 129. Nott, D. J., L. Marshall, and J. Brown (2012), Generalized likelihood uncertainty estimation  
2084 (GLUE) and approximate Bayesian computation: What's the connection?, *Water Resources*  
2085 *Research*, 48(12), doi:10.1029/2011WR011128.
- 2086 130. Pachepsky, Y. A., G. Martinez, F. Pan, T. Wagener, and T. Nicholson (2016), Evaluating  
2087 hydrological model performance using information theory-based metrics, *Hydrology and Earth*  
2088 *System Sciences Discussions*, 2016, 1–24, doi:10.5194/hess-2016-46.

- 2089 131. Palmer, T. N. (2002), The economic value of ensemble forecasts as a tool for risk assessment:  
2090      From days to decades, *Quarterly Journal of the Royal Meteorological Society*, 128(581), 747–774,  
2091      doi:10.1256/0035900021643593.
- 2092 132. Pool, S., M. Vis, and J. Seibert (2018), Evaluating model performance: towards a non-parametric  
2093      variant of the kling-gupta efficiency, *Hydrological Sciences Journal*, 63(13-14), 1941–1953, doi:  
2094      10.1080/02626667.2018.1552002.
- 2095 133. Raftery, A. E., T. Gneiting, F. Balabdaoui, and M. Polakowski (2005), Using Bayesian model  
2096      averaging to calibrate forecast ensembles, *Monthly Weather Review*, 133(5), 1155–1174, doi:  
2097      10.1175/MWR2906.1.
- 2098 134. Ramsey, F. P. (1926), Truth and probability, in *The Foundations of Mathematics and other Logical  
2099      Essays*, edited by R. B. Braithwaite, chap. 7, pp. 156–198, McMaster University Archive for the  
2100      History of Economic Thought.
- 2101 135. Rathinasamy, M., R. Khosa, J. Adamowski, S. ch, G. Partheepan, J. Anand, and B. Narsimlu  
2102      (2014), Wavelet-based multiscale performance analysis: An approach to assess and improve  
2103      hydrological models, *Water Resources Research*, 50(12), 9721–9737, doi:10.1002/2013WR014650.
- 2104 136. Renard, B., D. Kavetski, G. Kuczera, M. Thyre, and S. W. Franks (2010), Understanding predictive  
2105      uncertainty in hydrologic modeling: The challenge of identifying input and structural errors, *Water  
2106      Resources Research*, 46(5), doi:10.1029/2009WR008328.
- 2107 137. Renard, B., D. Kavetski, E. Leblois, M. Thyre, G. Kuczera, and S. W. Franks (2011), Toward  
2108      a reliable decomposition of predictive uncertainty in hydrological modeling: Characterizing  
2109      rainfall errors using conditional simulation, *Water Resources Research*, 47(11), doi:10.1029/  
2110      2011WR010643.
- 2111 138. Reusser, D. E., T. Blume, B. Schaeffli, and E. Zehe (2009), Analysing the temporal dynamics  
2112      of model performance for hydrological models, *Hydrology and Earth System Sciences*, 13(7),  
2113      999–1018, doi:10.5194/hess-13-999-2009.
- 2114 139. Rockafellar, R. T. (1970), *Convex Analysis*, Princeton Mathematical Series, 472 pp., Princeton  
2115      University Press, Princeton, NJ.

- 2116 140. Roques, C., D. E. Rupp, and J. S. Selker (2017), Improved streamflow recession parameter  
2117 estimation with attention to calculation of  $- dQ/dt$ , *Advances in Water Resources*, 108, 29–43,  
2118 doi:10.1016/j.advwatres.2017.07.013.
- 2119 141. Rosenblatt, M. (1952), Remarks on a multivariate transformation, *The Annals of Mathematical  
2120 Statistics*, 23(3), 470–472.
- 2121 142. Roulston, M. S., and L. A. Smith (2002), Evaluating probabilistic forecasts using information  
2122 theory, *Monthly Weather Review*, 130(6), 1653 – 1660, doi:10.1175/1520-0493(2002)130<1653:  
2123 EPFUIT>2.0.CO;2.
- 2124 143. Rupp, D. E., and J. S. Selker (2006), Information, artifacts, and noise in  $dq/dt - q$  recession  
2125 analysis, *Advances in Water Resources*, 29(2), 154–160, doi:10.1016/j.advwatres.2005.03.019,  
2126 experimental Hydrology: A Bright Future.
- 2127 144. Sadegh, M., and J. A. Vrugt (2013), Bridging the gap between GLUE and formal statistical  
2128 approaches: approximate Bayesian computation, *Hydrology and Earth System Sciences*, 17(12),  
2129 4831–4850, doi:10.5194/hess-17-4831-2013.
- 2130 145. Sadegh, M., J. Vrugt, H. Gupta, and C. Xu (2016), The soil water characteristic as new class  
2131 of closed-form parametric expressions for the flow duration curve, *Journal of Hydrology*, 535,  
2132 438–456, doi:10.1016/j.jhydrol.2016.01.027.
- 2133 146. Savage, L. J. (1971), Elicitation of personal probabilities and expectations, *Journal of the American  
2134 Statistical Association*, 66(336), 783–801, doi:doi.org/10.2307/2284229.
- 2135 147. Scharnagl, B., J. A. Vrugt, H. Vereecken, and M. Herbst (2010), Information content of incubation  
2136 experiments for inverse estimation of pools in the rothamsted carbon model: a bayesian perspective,  
2137 *Biogeosciences*, 7(2), 763–776, doi:10.5194/bg-7-763-2010.
- 2138 148. Scharnagl, B., S. C. Iden, W. Durner, H. Vereecken, and M. Herbst (2015), Inverse modelling of  
2139 in situ soil water dynamics: accounting for heteroscedastic, autocorrelated, and non-Gaussian  
2140 distributed residuals, *Hydrology and Earth System Sciences Discussions*, 12, 2155–2199.

- 2141 149. Scheuerer, M., and T. M. Hamill (2015), Variogram-based proper scoring rules for probabilistic  
2142 forecasts of multivariate quantities, *Monthly Weather Review*, 143(4), 1321–1334, doi:10.1175/  
2143 mwr-d-14-00269.1.
- 2144 150. Scheuerer, M., and D. Möller (2015), Probabilistic wind speed forecasting on a grid based on  
2145 ensemble model output statistics, *The Annals of Applied Statistics*, 9(3), 1328–1349.
- 2146 151. Schöniger, A., T. Wöhling, L. Samaniego, and W. Nowak (2014), Model selection on solid ground:  
2147 Rigorous comparison of nine ways to evaluate bayesian model evidence, *Water Resources Research*,  
2148 50(12), 9484–9513, doi:10.1002/2014WR016062.
- 2149 152. Schoups, G., and J. A. Vrugt (2010), A formal likelihood function for parameter and predictive  
2150 inference of hydrologic models with correlated, heteroscedastic, and non-Gaussian errors, *Water  
2151 Resources Research*, 46(10), doi:10.1029/2009WR008933.
- 2152 153. Schoups, G., J. A. Vrugt, F. Fenicia, and N. C. van de Giesen (2010), Corruption of accuracy and  
2153 efficiency of Markov chain Monte Carlo simulation by inaccurate numerical implementation of  
2154 conceptual hydrologic models, *Water Resources Research*, 46(10), doi:10.1029/2009WR008648.
- 2155 154. Schwemmle, R., D. Demand, and M. Weiler (2021), Technical note: Diagnostic efficiency –  
2156 specific evaluation of model performance, *Hydrology and Earth System Sciences*, 25, 2187–2198,  
2157 doi:10.5194/hess-25-2187-2021.
- 2158 155. Shamir, E., B. Imam, E. Morin, H. V. Gupta, and S. Sorooshian (2005), The role of hydrograph  
2159 indices in parameter estimation of rainfall–runoff models, *Hydrological Processes*, 19(11), 2187–  
2160 2207, doi:10.1002/hyp.5676.
- 2161 156. Shannon, C. E. (1948a), A mathematical theory of communication, *Bell System Technical Journal*,  
2162 27(3), 379–423, doi:10.1002/j.1538-7305.1948.tb01338.x.
- 2163 157. Shannon, C. E. (1948b), A mathematical theory of communication, *Bell System Technical Journal*,  
2164 27(4), 623–656, doi:10.1002/j.1538-7305.1948.tb00917.x.
- 2165 158. Shuford, E. H., A. Albert, and H. E. Massengill (1966), Admissible probability measurement  
2166 procedures, *Psychometrika*, 31, 125–145.

- 2167 159. Smirnov, N. (1948), Table for estimating the goodness of fit of empirical distributions, *The Annals*  
2168      *of Mathematical Statistics*, 19(2), 279–281, doi:10.1214/aoms/1177730256.
- 2169 160. Smolyak, S. A. (1963), Quadrature and interpolation formulas for tensor products of certain  
2170      classes of functions, *Doklady Akademii Nauk SSSR*, 148, 1042–1045.
- 2171 161. Sobol', I. M., and B. V. Shukhman (1995), Integration with quasirandom sequences: Nu-  
2172      mercial experience, *International Journal of Modern Physics C*, 6(2), 263–275, doi:10.1142/  
2173      S0129183195000204.
- 2174 162. Sorooshian, S., and J. A. Dracup (1980), Stochastic parameter estimation procedures for hydrologic  
2175      rainfall-runoff models: Correlated and heteroscedastic error cases, *Water Resources Research*,  
2176      16(2), 430–442, doi:10.1029/WR016i002p00430.
- 2177 163. Spear, R. C., and G. Hornberger (1980), Eutrophication in peel inlet-ii. identification of critical  
2178      uncertainties via generalized sensitivity analysis, *Water Research*, 14(1), 43–49.
- 2179 164. Spear, R. C., Q. Cheng, and S. L. Wu (2020), An example of augmenting regional sensitivity  
2180      analysis using machine learning software, *Water Resources Research*, 56(4), 1–16, doi:10.1029/  
2181      2019wr026379.
- 2182 165. Staël von Holstein, C.-A. S. (1970), A family of strictly proper scoring rules which are sensitive to  
2183      distance, *Journal of Applied Meteorology*, 9, 360–364.
- 2184 166. Stedinger, J. R., and G. D. Tasker (1985), Regional hydrologic analysis: 1. ordinary, weighted,  
2185      and generalized least squares compared, *Water Resources Research*, 21(9), 1421–1432, doi:  
2186      10.1029/wr021i009p01421.
- 2187 167. Storch, H. v., and F. W. Zwiers (1999), *Statistical Analysis in Climate Research*, p. 484, 1 ed.,  
2188      Cambridge University Press, doi:10.1017/CBO9780511612336.
- 2189 168. Székely, G. J. (2003),  $\mathcal{E}$ -statistics: The energy of statistical samples, *Tech. rep.*, Department of  
2190      Mathematics and Statistics, Bowling Green State University, Bowling Green, OH, USA.

- 2191 169. Tashie, A., T. Pavelsky, and L. E. Band (2020), An empirical reevaluation of streamflow recession  
2192 analysis at the continental scale, *Water Resources Research*, 56(1), e2019WR025,448, doi:10.1029/  
2193 2019WR025448.
- 2194 170. Tasker, G. D. (1980), Hydrologic regression with weighted least squares, *Water Resources Research*,  
2195 16(6), 1107–1113, doi:10.1029/wr016i006p01107.
- 2196 171. Tegos, S. A., G. K. Karagiannidis, P. D. Diamantoulakis, and N. D. Chatzidiamantis (2022), New  
2197 results for pearson type iii family of distributions and application in wireless power transfer, *IEEE*  
2198 *Internet of Things Journal*, 9(23), 24,038–24,050, doi:10.1109/JIOT.2022.3189220.
- 2199 172. Thielen, J., J. Schaake, R. Hartman, and R. Buizza (2008), Aims, challenges and progress of the  
2200 hydrological ensemble prediction experiment (HEPEX) following the third HEPEX workshop held  
2201 in Stresa 27 to 29 june 2007, *Atmospheric Science Letters*, 9, 29–35, doi:10.1002/asl.168.
- 2202 173. Thomas, B. F., R. M. Vogel, C. N. Kroll, and J. S. Famiglietti (2013), Estimation of the base  
2203 flow recession constant under human interference, *Water Resources Research*, 49(11), 7366–7379,  
2204 doi:10.1002/wrcr.20532.
- 2205 174. Thorarinsdottir, T. L., T. Gneiting, and N. Gissibl (2013), Using proper divergence functions  
2206 to evaluate climate models, *SIAM/ASA Journal on Uncertainty Quantification*, 1(1), 522–534,  
2207 doi:10.1137/130907550.
- 2208 175. Thyer, M., B. Renard, D. Kavetski, G. Kuczera, S. W. Franks, and S. Srikanthan (2009), Critical  
2209 evaluation of parameter consistency and predictive uncertainty in hydrological modeling: A  
2210 case study using Bayesian total error analysis, *Water Resources Research*, 45(12), doi:10.1029/  
2211 2008WR006825.
- 2212 176. Tsyplakov, A. (2011), Evaluating density forecasts: A comment, *SSRN*, doi:10.2139/ssrn.1907799.
- 2213 177. Unger, D. A. (1985), A method to estimate the continuous ranked probability score, in *Preprints of*  
2214 *the Ninth Conference on Probability and Statistics in Atmospheric Sciences*, pp. 206–213, American  
2215 Meteorological Society, Virginia Beach, Virginia, USA.

- 2216 178. Villez, K. (2017), Analytical expressions to compute the Continuous Ranked Probability Score  
2217 (CRPS), *Tech. Rep.* 4, Eawag, Aquatic Research, Swiss Federal Institute of Aquatic Science and  
2218 Technology, Dübendorf, Switzerland.
- 2219 179. Vogel, R. M., and N. M. Fennessey (1994), Flow-duration curves. i: New interpretation and  
2220 confidence intervals, *Journal of Water Resources Planning and Management*, 120(4), 485–504,  
2221 doi:10.1061/(ASCE)0733-9496(1994)120:4(485).
- 2222 180. Volpi, E., G. Schoups, G. Firmani, and J. A. Vrugt (2017), Sworn testimony of the model evidence:  
2223 Gaussian mixture importance (GAME) sampling, *Water Resources Research*, 53(7), 6133–6158,  
2224 doi:10.1002/2016WR020167.
- 2225 181. Vrugt, J. A. (2016), Markov chain Monte Carlo simulation using the DREAM software package:  
2226 Theory, concepts, and matlab implementation, *Environmental Modelling & Software*, 75, 273–316,  
2227 doi:10.1016/j.envsoft.2015.08.013.
- 2228 182. Vrugt, J. A. (2018), MODELAVG: A MATLAB toolbox for postprocessing of model ensembles,  
2229 *Tech. rep.*, University of California, Irvine.
- 2230 183. Vrugt, J. A., and K. J. Beven (2018), Embracing equifinality with efficiency: Limits of acceptability  
2231 sampling using the DREAM<sub>(LOA)</sub> algorithm, *Journal of Hydrology*, 559, 954–971, doi:10.1016/j.jhydrol.  
2232 2018.02.026.
- 2233 184. Vrugt, J. A., and D. Y. de Oliveira (2022), Confidence intervals of the kling-gupta efficiency,  
2234 *Journal of Hydrology*, 612, 127,968, doi:10.1016/j.jhydrol.2022.127968.
- 2235 185. Vrugt, J. A., and B. A. Robinson (2007), Treatment of uncertainty using ensemble methods:  
2236 Comparison of sequential data assimilation and Bayesian model averaging, *Water Resources  
2237 Research*, 43(1), doi:10.1029/2005WR004838.
- 2238 186. Vrugt, J. A., and M. Sadegh (2013a), Toward diagnostic model calibration and evaluation:  
2239 Approximate Bayesian computation, *Water Resources Research*, 49(7), 4335–4345, doi:10.1002/  
2240 wrcr.20354.

- 2241 187. Vrugt, J. A., W. Bouten, H. V. Gupta, and S. Sorooshian (2002), Toward improved identifiability  
2242 of hydrologic model parameters: The information content of experimental data, *Water Resources*  
2243 *Research*, 38(12), 48–1–48–13, doi:10.1029/2001WR001118.
- 2244 188. Vrugt, J. A., H. V. Gupta, L. A. Bastidas, W. Bouten, and S. Sorooshian (2003), Effective and  
2245 efficient algorithm for multiobjective optimization of hydrologic models, *Water Resources Research*,  
2246 39(8), doi:10.1029/2002WR001746.
- 2247 189. Vrugt, J. A., M. P. Clark, C. G. H. Diks, Q. Duan, and B. A. Robinson (2006), Multi-objective  
2248 calibration of forecast ensembles using Bayesian model averaging, *Geophysical Research Letters*,  
2249 33(19), doi:10.1029/2006GL027126.
- 2250 190. Vrugt, J. A., C. G. H. Diks, and M. P. Clark (2008), Ensemble Bayesian model averaging  
2251 using Markov chain Monte Carlo sampling, *Environmental Fluid Mechanics*, 8(5), 579–595,  
2252 doi:10.1007/s10652-008-9106-3.
- 2253 191. Vrugt, J. A., D. Yumi de Oliveria, G. Schoups, and C. G. H. Diks (2022), On the use of distribution-  
2254 adaptive likelihood functions: Generalized and universal likelihood functions, scoring rules and  
2255 multi-criteria ranking, *Journal of Hydrology*, 615, 128,542, doi:10.1016/j.jhydrol.2022.128542.
- 2256 192. Wagener, T., N. McIntyre, M. J. Lees, H. S. Wheater, and H. V. Gupta (2003), Towards reduced  
2257 uncertainty in conceptual rainfall-runoff modelling: dynamic identifiability analysis, *Hydrological*  
2258 *Processes*, 17(2), 455–476, doi:10.1002/hyp.1135.
- 2259 193. Weijs, S. V., G. Schoups, and N. van de Giesen (2010a), Why hydrological predictions should be  
2260 evaluated using information theory, *Hydrology and Earth System Sciences*, 14(12), 2545–2558,  
2261 doi:10.5194/hess-14-2545-2010.
- 2262 194. Weijs, S. V., R. van Nooijen, and N. van de Giesen (2010b), Kullback-leibler divergence as a  
2263 forecast skill score with classic reliability-resolution-uncertainty decomposition, *Monthly Weather*  
2264 *Review*, 138(9), 3387–3399, doi:10.1175/2010MWR3229.1.
- 2265 195. Welles, E., S. Sorooshian, G. Carter, and B. Olsen (2007), Hydrologic verification: A call for  
2266 action and collaboration, *Bulletin of the American Meteorological Society*, 88(4), 503–512, doi:  
2267 10.1175/BAMS-88-4-503.

- 2268 196. Westerberg, I. K., J.-L. Guerrero, P. M. Younger, K. J. Beven, J. Seibert, S. Halldin, J. E. Freer,  
2269 and C.-Y. Xu (2011), Calibration of hydrological models using flow-duration curves, *Hydrology*  
2270 and *Earth System Sciences*, 15(7), 2205–2227, doi:10.5194/hess-15-2205-2011.
- 2271 197. Winkler, R. L., et al. (1996), Scoring rules and the evaluation of probabilities, *Test*, 5(1), 1–60,  
2272 doi:10.1007/BF02562681.
- 2273 198. Ye, M., P. D. Meyer, and S. P. Neuman (2008), On model selection criteria in multimodel analysis,  
2274 *Water Resources Research*, 44(3), doi:10.1029/2008WR006803.
- 2275 199. Yilmaz, K. K., H. V. Gupta, and T. Wagener (2008), A process-based diagnostic approach to model  
2276 evaluation: Application to the NWS distributed hydrologic model, *Water Resources Research*,  
2277 44(9), doi:10.1029/2007WR006716.
- 2278 200. Zheng, J., and H. You (2013), A new model-independent method for change detection in multi-  
2279 temporal SAR images based on radon transform and Jeffrey divergence, *IEEE Geoscience and*  
2280 *Remote Sensing Letters*, 10(1), 91–95, doi:10.1109/LGRS.2012.2193659.

## Figure Captions

2281

2282 **Figure 1:** Probability density function of the uniform true distribution  $Q = \mathcal{U}(a, b)$  and symmetric  
 2283 triangular forecast distribution  $P = \mathcal{T}(a, b)$  of quantity  $\Omega = [a, b]$ .

2284 **Figure 2:** Illustration of the computation of the relative entropy  $d_{\text{KL}}(Q, P)$  on a sample space  $\Omega =$   
 2285  $\{0, 1, 2, 3, 4\}$ : (a) *true distribution*,  $Q = \mathcal{B}(n, p)$ , and (b) *forecast distribution*,  $P = \mathcal{U}_d(n)$ . According to  
 2286 data,  $\Omega$  follows a binomial distribution with  $n = 4$ ,  $p = \frac{1}{2}$  and PMF,  $f_{\mathcal{B}}(\omega, n, p) = c(n, \omega)p^{\omega}(1-p)^{n-\omega}$ ,  
 2287 where  $c(a, b) = a!/b!(a - b)!$  denotes the binomial coefficient and ! is the factorial function. Theory  
 2288 predicts a discrete uniform distribution for  $\omega$  with equal density,  $f_{\mathcal{U}_d}(\omega, n) = 1/n$ , for all  $n = 5$  values.

2289 **Figure 3:** (a) Illustration of the probability integral transform and (b) interpretation of the so-called  
 2290 quantile-quantile plot (adapted from *Laio and Tamea (2007)*).

2291 **Figure 4:** Hypothetical *true distribution*  $Q = \mathcal{G}(3.36, 0.64)$  (gray) and different *probabilistic forecasts*  $P$   
 2292 (in color) with an equal (a) coefficient of variation,  $C_v = 0.546$ , and (b) width,  $W = 4.5$ , at  $\alpha = 0.05$   
 2293 using lower and upper quantiles of 0.5 and 5.0, respectively. The peak of the *true distribution*  
 2294 is equal to the verifying observation,  $\omega$ . The distribution forecasts of the left graph are used in later  
 2295 studies:  $P_1 = \mathcal{U}(0.11, 3.89)$ ,  $P_2 = \mathcal{N}(2.26, 1.52)$ ,  $P_3 = \mathcal{GEV}(0.04, 1, 1.86)$ ;  $P_4 = \mathcal{LN}(0.80, 0.26)$  and  
 2296  $P_5 = \mathcal{GP}(-0.28, 2, 0.73)$ .

2297 **Figure 5:** Generalized entropy function  $H(p) = p(p - 1)$  (blue curve) of the quadratic score for a  
 2298 dichotomous event  $\Omega = \{1, 0\}$  with *probability forecast*  $(p, 1 - p)$  and *true probability*  $(q, 1 - q)$  with  
 2299  $p, q \in [0, 1]$ . We present the values of the quadratic scoring rule  $\mathcal{S}(p, q)$  at  $p$  and  $q$  (solid black dots) and  
 2300 display the so-called Bregman divergence,  $d(p, q)$ . For any probability forecast,  $p \in [0, 1]$ , the expected  
 2301 score,  $\mathcal{S}(p, q) = qS(p, 1) + (1-q)S(p, 0)$ , equals the ordinate of the tangent to  $H$  at  $p$  (solid gray line) when  
 2302 evaluated at  $q \in [0, 1]$ . In particular, the scores,  $S(p, 0) = H(p) - pH'(p)$  and  $S(p, 1) = H(p) + (1-p)H'(p)$ ,  
 2303 equal the tangent at  $q = 0$  and  $q = 1$ , respectively. The divergence,  $d_{\text{QS}}(p, q) = \mathcal{S}(q, q) - \mathcal{S}(p, q)$ , is  
 2304 equal to the difference between  $H(q)$  and the tangent at  $p$  when evaluated at  $q$  (Adapted after Fig. 1 of  
 2305 *Gneiting and Raftery (2007)* and Fig. 8 of *Buja et al. (2005)*).

2306 **Figure 6:** Binary event,  $\Omega = \{1, 0\}$ : Expected value of the quadratic (green), logarithmic (red) and  
 2307 spherical (blue) scoring rules as function of the *true probability*  $q = q_1$  of the first event, (a)  $q = 0$ , (b)  
 2308  $q = 0.1$ , (c)  $q = 0.2$ , (d)  $q = 0.3$ , (e)  $q = 0.4$ , (f)  $q = 0.5$ , (g)  $q = 0.6$ , (h)  $q = 0.7$ , (i)  $q = 0.8$ , and (j)

2309  $q = 0.9$ .

2310 **Figure 7:** Illustration of the standard normal Lebesgue density on sample space  $\Omega = [-3, 3]$  with range of  
2311  $\mathcal{N}(0, 1)$  partitioned into 14 intervals. The Lebesgue measure  $\mu(\omega_k)$  is equal to the length of the interval  
2312 ( $=$  bin width) containing  $\omega_k$ . The Lebesgue density  $f(\omega)$  is constant in each bin  $\omega_1, \dots, \omega_m$  of  $\omega$  values,  
2313 where  $m = 29$ . The sum of the areas of the rectangles is equal to the Lebesgue integral.

2314 **Figure 8:** Score divergence  $d_{\text{LinS}}(P, Q)$  of the linear scoring rule  $S_{\text{LinS}}(P, \omega) = f_P(\omega)$  for a uniform  
2315 probabilistic forecast  $P$  under a standard Gaussian true distribution  $Q$  and symmetric interval  $[-\varepsilon, \varepsilon]$   
2316 with  $\varepsilon \in (0, 3]$ .

2317 **Figure 9:** Graphical explanation of the continuous ranked probability scoring rule for a hypothetical  
2318 streamflow forecast CDF ( $F_P$ , black line) and verifying observation ( $\omega$ , red dot).

2319 **Figure 10:** Hypothetical distribution forecast  $P$  of the discharge  $\Omega$  (mm/d) and traces of the quadratic  
2320 (green), logarithmic (red), spherical (blue) and continuous ranked probability (yellow) scoring rules across  
2321 the streamflow distribution computed using their numerical definitions in Table 6. The probabilistic  
2322 forecast is a gamma distribution  $P = \mathcal{G}(a, b)$  with shape and scale parameters  $a = 3$  (-) and  $b = 1$  mm/d,  
2323 respectively, and PDF  $f_g(\omega, a, b) = \Gamma^{-1}(a)b^{-a}\omega^{a-1} \exp(-\omega/b)$ , where  $\Gamma(x)$  is the gamma function.

2324 **Figure 11:** Streamflow predictions of the ABC (3), GR4J (4), HYMOD (5), TOPMO (6), AWBM (8),  
2325 NAM (9), HBV (9) and SAC-SMA (13) conceptual watershed models for a small period of the  $n = 3,000$   
2326 day training data record. The number between parenthesis lists the number of calibration parameters.  
2327 The solid red circles correspond to the daily measured discharges.

2328 **Figure 12:** Weighted average BMA forecast (solid blue line) and 50% (dim gray), 75% (medium gray)  
2329 and 95% (light gray) quantiles of the BMA forecast distribution for a representative 270-day period of  
2330 the training data set using the (a) normal, (b) lognormal, (c) generalized normal, (d) gamma and (e)  
2331 Weibull distributions with a model-dependent (= sole) forecast variance. The red circles are the daily  
2332 discharge observations.

2333 **Figure 13:** (a) 50%, 75% and 95% prediction intervals of the BMA forecast density for the lognormal  
2334 distribution of Table 8 with a constant model-dependent (= sole) forecast variance and corresponding  
2335 traces of the (b) quadratic, (c) logarithmic, (d) spherical, (e) continuous ranked probability and (f)  
2336 energy scoring rules.

2337 **Figure 14:** Relationship between the ensemble size,  $K$ , and the time-averaged value of the logarithmic  
2338 scoring rule for HYMOD (red line), Hmodel (blue line) and SAC-SMA (green line) using daily discharge  
2339 data from the Leaf River watershed. The dashed vertical lines point out the optimum size of the ensemble.

2340 **Figure 15:** The 50, 75, 95 and 100% simulation uncertainty intervals of the (a) HYMOD, (b) Hmodel  
2341 and (c) SAC-SMA conceptual hydrologic models for the optimum ensemble according to the logarithmic  
2342 scoring rule.

2343 **Figure 16:** Preliminary results of the cumulative ranked exceedance probability divergence score: (a)  
2344 measured FDC (red dots) and SAC-SMA model simulated reliability functions (gray lines) of the  $K = 50$   
2345 discharge records of the ensemble with lowest values of the CREPS divergence and (b,c) bivariate scatter  
2346 plots of the CREPS divergence and (b) mean taxicab distance,  $\bar{d}_T(P, Q)$  and (c) normal log-likelihood  
2347  $\ell(\boldsymbol{\theta}|\boldsymbol{\omega})$  of Equation (79) using all discharge records of the GLUE ensemble.

2348 **Figure 17:** Traces of the expected value of the interval score  $S_{IS}^\alpha(P, \omega)$  as function of the lower endpoint  $l$   
2349 of the  $100(1 - \alpha)\%$  prediction interval using  $\alpha = 0.05$  (red),  $\alpha = 0.25$  (blue) and  $\alpha = 0.50$  (green). The  
2350 colored dots are a projection of the maximum interval score on the  $x$ -axis.

2351 **Figure 18:** Scatter plots (blue squares) of the  $\log_{10}(y)$  versus  $\log_{10}(\text{dy}/\text{dt})$  relationship for the different  
2352 models of the BMA discharge ensemble. The red dots correspond to the measured discharge record.

## Table Captions

2354 **Table 1:** Time-averaged performance measures,  $\bar{M}(\mathbf{P}, \boldsymbol{\omega})$ , of distribution forecasts  $\mathbf{P} = \{P_1, \dots, P_n\}$  and  
 2355 verifying observations  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_n)^\top$ .

2356 **Table 2:** Quadratic, logarithmic and pseudospherical scoring rules for categorical variables: Entropy  
 2357 function, scoring rule, expected score function and divergence function for a *distribution forecast*  $\mathbf{p} =$   
 2358  $(p_1, \dots, p_m)^\top$  on the convex class  $\mathcal{P} = \mathcal{P}_m$  of  $m \geq 2$  mutually exclusive and collectively exhaustive  
 2359 events,  $\Omega = \{1, \dots, m\}$ . The  $m$ -vector  $\mathbf{q} = (q_1, \dots, q_m)^\top$  lists the true event probabilities.

2360 **Table 3:** Strictly proper scoring rules for a dichotomous event (*rain* and *no rain*) with *probability forecast*  
 2361  $\mathbf{p} = (p, 1 - p)$  on  $\Omega = \{1, 0\}$  with  $p \in [0, 1]$ .

2362 **Table 4:** Entropy,  $H(\mathbf{p})$ , expectation,  $\mathcal{S}(\mathbf{p}, \mathbf{q})$  and divergence,  $d(\mathbf{p}, \mathbf{q})$  of the *strictly proper* categorical  
 2363 scoring rules of Table 2 for distribution forecasts  $P_1, \dots, P_5$  displayed in Fig. 4a using  $m = 60$  discrete  
 2364 values,  $\Omega = \frac{1}{20}\{1, 3, 5, \dots, 117, 119\}$ . Column  $R$  lists the rank of the distribution forecasts. The bottom  
 2365 row presents the values for a perfect distribution forecast,  $P = Q$ .

2366 **Table 5:** Entropy,  $H(\mathbf{p})$ , expectation,  $\mathcal{S}(\mathbf{p}, \mathbf{q})$  and divergence,  $d(\mathbf{p}, \mathbf{q})$  of the *strictly proper* categorical  
 2367 scoring rules of Table 2 for the *true* and *forecasted* rainfall probabilities of Table E.1.

2368 **Table 6:** Summary of *strictly proper* scoring rules for a density forecast  $f_P$  and verifying observation  $\omega$ .  
 2369 The numerical form assumes that the forecast distribution  $P$  is a  $m$ -member ensemble  $(y_1, \dots, y_m)^\top$ .

2370 **Table 7:** Mean values of the quadratic, logarithmic, spherical and continuous ranked probability scoring  
 2371 rules for the distribution forecasts  $P_1, \dots, P_5$  portrayed in Fig. 4a. The last column reports the mean  
 2372 scores for a perfect distribution forecast,  $P = Q$ .

2373 **Table 8:** Formulation, coefficients and unknown parameters of the 1. generalized normal, 2. lognormal, 3.  
 2374 gamma and 4. Weibull predictive PDFs used in the BMA forecast density of Equation (70).

2375 **Table 9:** Time-averaged values of the *strictly proper* scoring rules of Table 6 for the BMA density forecast  
 2376  $f_{P_t}(y|\boldsymbol{\beta}, \boldsymbol{\psi})$ ;  $t = (1, \dots, n)$  of Equation (70) using the normal ( $\tau = 2$ ), lognormal, generalized normal,  
 2377 gamma and Weibull predictive PDFs with a group or model-dependent (= sole) constant variance. We also  
 2378 list the performance metrics,  $R_l$ ,  $C_v$ ,  $C$  and  $W$  of Table 1 and report the log-likelihood,  $\ell(\boldsymbol{\beta}, \boldsymbol{\psi}|\boldsymbol{\omega})$ , Root  
 2379 Mean Square Error (RMSE), coefficient of determination ( $R^2$ ) and KG efficiency of the weighted-average  
 2380 BMA forecast. The bottom row lists the number  $d$  of unknown BMA parameters.

2381 **Table 10:** Pearson correlation coefficients of the log-likelihood  $\ell(\boldsymbol{\beta}, \boldsymbol{\psi} | \boldsymbol{\omega})$  of the weighted-average BMA  
 2382 forecast and time-averaged values of the QS, LS, SS, CRPS and ES *strictly proper* scoring rules, improper  
 2383 performance metrics of Table 1 and RMSE,  $R^2$  and KGE scoring functions.

2384 **Table 11:** Time-averaged values of the moments (skew and kurtosis) and performance metrics ( $R_l$ ,  $C_v$ ,  
 2385  $C$  and  $W$ ) of the distribution forecasts simulated by HYMOD, Hmodel and the SAC-SMA model and  
 2386 RMSE,  $R^2$  and KGE scoring functions of the mean functional of the forecasts. The coverage  $C$  of the  
 2387  $\gamma = 1 - \alpha$  prediction intervals corresponds to the significance levels  $\alpha = 0.5$ , 0.75 and 0.95 of Fig. 15.  
 2388 The second column  $d$  lists the number of unknown parameters.

2389 **Table 12:** Unconditional,  $\bar{p}$ , and conditional,  $\pi$ , probabilities of the watershed models estimated from  
 2390 the 3,000-day training data record:  $\pi_{jk} = \mathbb{P}(\omega = y_j | y_k)$  is the probability of  $y_j$  given that  $y_k$  is the best  
 2391 forecast in the ensemble at the previous time.

2392 **Table 13:** Time-averaged expected value of the *strictly proper* quadratic scoring rule (see Table 6),  
 2393 entropy, resolution and reliability for the BMA density forecasts of Equation (70) using the normal,  
 2394 lognormal, generalized normal, gamma and Weibull distributions with a constant group and sole variance,  
 2395 respectively. The bottom row completes the decomposition of Equation (100).

2396 **Table 14:** Mean values of the Bayes factors  $B_{i,j}$  of the distribution forecasts  $P_1, \dots, P_5$  displayed in Fig.  
 2397 4a from application of Equation (111) to the sum of the logarithmic scores  $S_{LS}(P, \omega)$  in units of bits  
 2398 ( $b = e$ ) of the  $n = 10^4$  outcomes  $\omega_1, \dots, \omega_n$ . The last row and column correspond to a perfect distribution  
 2399 forecast,  $P = Q$ . *Kass and Raftery* (1995) (P. 777) categorize the Bayes factors and present descriptive  
 2400 statements on the strength of the evidence;  $B_{i,j} = 0$  favors in strongest possible terms the null hypothesis,  
 2401  $P_j$ , and vice-versa  $B_{i,j} > 75$  means that there is decisive evidence against  $P_j$ .

## 2402 Appendix A: On Gibbs' inequality

2403 Gibbs' inequality

$$2404 \quad \mathbb{H}(Q, P) \geq \mathbb{H}(Q) \quad (\text{A.1})$$

2405 was presented by the American scientist Josiah Willard Gibbs (1839-1903) and states that the cross-  
 2406 entropy  $\mathbb{H}(Q, P)$  of two probability distributions  $Q$  and  $P$  will always exceed the entropy  $\mathbb{H}(Q) = \mathbb{H}(Q, Q)$   
 2407 of distribution  $Q$  alone unless  $P = Q$  then  $\mathbb{H}(Q, P) = \mathbb{H}(Q)$ . Different mathematical proofs exist of this  
 2408 inequality. For completeness, we present one of them in this Appendix.

2409 Suppose  $Q$  and  $P$  are discrete probability distributions on a common sample space  $\Omega$ . If  $x$  is a possible  
 2410 outcome then  $q(x) \geq 0$  and  $p(x) \geq 0$  denote the probability for  $x \in \Omega$  with  $\sum_{x \in \Omega} q(x) = 1$  and  
 2411  $\sum_{x \in \Omega} p(x) = 1$ . Equation (A.1) may now be written in discretized form

$$2412 \quad \mathbb{H}(\mathbf{q}, \mathbf{p}) \geq \mathbb{H}(\mathbf{q}) \quad (\text{A.2})$$

2413 where  $\mathbf{q}$  and  $\mathbf{p}$  are vectors with probabilities of  $Q$  and  $P$  for all events of  $\Omega$ . Now we can write

$$2414 \quad \sum_{x \in \Omega} q(x) \log_b \left( \frac{1}{p(x)} \right) \geq \sum_{x \in \Omega} q(x) \log_b \left( \frac{1}{q(x)} \right), \quad (\text{A.3})$$

2416 which is equal to

$$2417 \quad \sum_{x \in \Omega} q(x) \log_b(q(x)) - \sum_{x \in \Omega} q(x) \log_b(p(x)) \geq 0, \quad (\text{A.4})$$

2419 and simplifies further to

$$2420 \quad \sum_{x \in \Omega} q(x) \left( \log_b(q(x)) - \log_b(p(x)) \right) \geq 0$$

$$2421 \quad \sum_{x \in \Omega} q(x) \log_b \left( \frac{q(x)}{p(x)} \right) \geq 0$$

$$2422 \quad \implies d_{\text{KL}}(Q, P) \geq 0. \quad (\text{A.5})$$

2424 Thus, to prove Gibbs' inequality we need to demonstrate that the relative entropy,  $d_{\text{KL}}(Q, P)$ , is  
 2425 nonnegative. We define the function  $t(x) = p(x)/q(x)$ . This function satisfies the following condition

$$2426 \quad \log_b(t(x)) \leq t(x) - 1 \quad (\text{A.6})$$

2427 for all  $t(x) > 0$  and  $b > 0$  with equality if and only if  $t(x) = 1$  and, thus,  $P = Q$ . The above inequality  
 2428 may be written as follows

$$2429 \quad -\log_b \left( \frac{q(x)}{p(x)} \right) \leq \frac{p(x)}{q(x)} - 1, \quad (\text{A.7})$$

2431 and, thus, we yield

$$\log_b \left( \frac{q(x)}{p(x)} \right) \geq 1 - \frac{p(x)}{q(x)}. \quad (\text{A.8})$$

<sup>2434</sup> We can multiply both sides of Equation (A.8) with  $q(x)$

$$\begin{aligned} \sum_{x \in \Omega} q(x) \log_b \left( \frac{q(x)}{p(x)} \right) &\geq \sum_{x \in \Omega} q(x) \left( 1 - \frac{p(x)}{q(x)} \right) \\ &\geq \sum_{x \in \Omega} q(x) - \sum_{x \in \Omega} p(x), \end{aligned} \tag{A.9}$$

2438 to arrive at the following inequality

$$d_{\text{KL}}(Q, P) \geq 1 - \sum_{x \in \Omega} p(x). \quad (\text{A.10})$$

<sup>2441</sup> Thus, for Gibbs' inequality to hold we simply need to show that

$$1 - \sum_{x \in \Omega} p(x) \geq 0 \iff \sum_{x \in \Omega} p(x) \leq 1. \quad (\text{A.11})$$

In plain words, the sum of the probabilities of distribution  $P$  at the collection of outcomes  $x \in \Omega$  cannot exceed unity. This condition will always be satisfied as, (i) all non-zero  $q(x)$  values will sum to one and so do their  $p(x)$  values, and (ii) all points  $x \in \Omega$  at which  $p(x) > 0$  but  $q(x) = 0$  do not contribute to the relative entropy as  $\lim_{q \downarrow 0} q \log(q) = 0$ . Hence, the sum of the  $p(x)$ 's at which  $q(x) > 0$  will almost surely be smaller than one unless  $P = Q$  then  $\mathbb{H}(Q, P) = \mathbb{H}(Q)$  and  $d_{\text{KL}}(Q, P) = 0$ .

## 2449 Appendix B: Relative Entropy Triangular and Uniform Distributions

2450 In this Appendix, we present the derivation of the relative entropy of the triangular forecast distribution  
 2451  $P$  and uniform true distribution  $Q$  displayed in Figure 1. If we enter the analytic expressions of the  
 2452 PDFs of  $P$  and  $Q$  into the integral of Equation (1) we yield

$$\begin{aligned} 2453 \quad d_{\text{KL}}(Q, P) &= \int_a^b \frac{1}{b-a} \log_b \left( \frac{(b-a)^2}{(b-a)(2(b-a)-2|a+b-2x|)} \right) dx \\ 2454 &= \frac{1}{b-a} \int_a^c \log_b \left( \frac{b-a}{2(b-a)-2(a+b-2x)} \right) dx + \frac{1}{b-a} \int_c^b \log_b \left( \frac{b-a}{2(b-a)-2(2x-a-b)} \right) dx \\ 2455 &= \frac{1}{b-a} \int_a^c \log_b \left( \frac{1}{4} \frac{b-a}{x-a} \right) dx + \frac{1}{b-a} \int_c^b \log_b \left( \frac{1}{4} \frac{b-a}{b-x} \right) dx \\ 2456 &= \frac{1}{b-a} \left| (x-a) \left[ \log_b \left( \frac{1}{4} \frac{a-b}{a-x} \right) + 1 \right] \right|_a^c + \frac{1}{b-a} \left| (x-b) \left[ \log_b \left( \frac{1}{4} \frac{b-a}{b-x} \right) + 1 \right] \right|_c^b. \\ 2457 \end{aligned} \quad (\text{B.1})$$

2458 At the midpoint  $c$ , we yield  $x-a = \frac{1}{2}(b-a)$  and  $x-b = \frac{1}{2}(a-b)$ , thus, the expression above simplifies to

$$\begin{aligned} 2459 \quad d_{\text{KL}}(Q, P) &= \left( \frac{\frac{1}{2}(b-a)}{b-a} \left[ \log_b \left( \frac{1}{4} \frac{a-b}{\frac{1}{2}(a-b)} \right) + 1 \right] - 0 \right) + \left( 0 - \frac{\frac{1}{2}(a-b)}{b-a} \left[ \log_b \left( \frac{1}{4} \frac{b-a}{\frac{1}{2}(b-a)} \right) + 1 \right] \right) \\ 2460 &= \frac{1}{2} \left[ \log_b \left( \frac{1}{2} \right) + 1 \right] + \frac{1}{2} \left[ \log_b \left( \frac{1}{2} \right) + 1 \right] \\ 2461 &= 1 - \log_b(2). \\ 2462 \\ 2463 \end{aligned} \quad (\text{B.2})$$

2464 We can swap the arguments  $Q$  and  $P$  and compute the relative entropy from  $Q$  to  $P$  to yield

$$\begin{aligned} 2465 \quad d_{\text{KL}}(P, Q) &= \int_a^b \frac{2(b-a)-2|a+b-2x|}{(b-a)^2} \log_b \left( \frac{(b-a)(2(b-a)-2|a+b-2x|)}{(b-a)^2} \right) dx \\ 2466 &= \int_a^c \frac{4(x-a)}{(b-a)^2} \log_b \left( \frac{4(x-a)}{b-a} \right) dx + \int_c^b \frac{4(b-x)}{(b-a)^2} \log_b \left( \frac{4(b-x)}{b-a} \right) dx \\ 2467 &= \left| \frac{(4a-4x)^2}{8(a-b)^2} \left[ \log_b \left( \frac{4(a-x)}{a-b} \right) - \frac{1}{2} \right] \right|_a^{\frac{a+b}{2}} + \left| -\frac{(4b-4x)^2}{8(a-b)^2} \left[ \log_b \left( \frac{4(x-b)}{a-b} \right) - \frac{1}{2} \right] \right|_a^{\frac{a+b}{2}} \\ 2468 &= \left( \frac{4(a-b)^2}{8(a-b)^2} \left[ \log_b \left( \frac{2(a-b)}{a-b} \right) - \frac{1}{2} \right] - 0 \right) + \left( 0 + \frac{4(b-a)^2}{8(a-b)^2} \left[ \log_b \left( \frac{2(a-b)}{a-b} \right) - \frac{1}{2} \right] - 0 \right) \\ 2469 &= \frac{1}{2} \left[ \log_b(2) - \frac{1}{2} \right] + \frac{1}{2} \left[ \log_b(2) - \frac{1}{2} \right] \\ 2470 &= \log_b(2) - \frac{1}{2}. \\ 2471 \end{aligned} \quad (\text{B.4})$$

2472 This concludes our derivation.

## 2473 Appendix C: KL divergence of two multivariate normal distributions

2474 Suppose that the *true* joint distribution  $Q$  of  $\mathbf{x} = (x_1, \dots, x_b)^\top$  and its probabilistic forecast,  $P$ , are each  
 2475 described by a multivariate normal distribution,  $\mathcal{N}_b(\boldsymbol{\mu}_Q, \boldsymbol{\Sigma}_Q)$  and  $\mathcal{N}_b(\boldsymbol{\mu}_P, \boldsymbol{\Sigma}_P)$ , respectively, with means,  
 2476  $\boldsymbol{\mu}_Q = [\mu_{Q,1} \dots \mu_{Q,b}]^\top$  and  $\boldsymbol{\mu}_P$ , and  $b \times b$  covariance matrices,  $\boldsymbol{\Sigma}_Q$  and  $\boldsymbol{\Sigma}_P$ , respectively. The probability  
 2477 density at  $\mathbf{x}$  is then equal to

$$2478 f_N(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{b/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right), \quad (\text{C.1})$$

2479 where  $|\cdot|$  is the determinant operator and the symbol  $\top$  denotes transpose. We can use the above  
 2480 expression to derive a closed-form expression for the KL-divergence,  $d_{\text{KL}}(Q, P)$ , of  $Q$  and  $P$  in  $\mathbb{R}^b$ . Indeed,  
 2481 we can write

$$2482 d_{\text{KL}}(Q, P) = \mathbb{E}_Q \left[ \log_e \left( \frac{Q(x)}{P(x)} \right) \right] \\ 2483 = \mathbb{E}_Q [\log_e(f_N(\mathbf{x}, \boldsymbol{\mu}_Q, \boldsymbol{\Sigma}_Q)) - \log_e(f_N(\mathbf{x}, \boldsymbol{\mu}_P, \boldsymbol{\Sigma}_P))]. \quad (\text{C.2})$$

To cancel the exponential function in Equation (C.1), the base of the logarithm must be fixed to Euler's number,  $e = 2.7182818\dots$ , and as a result the relative entropy,  $d_{\text{KL}}(Q, P)$ , has units of nats. If we admit  
 2485 to Equation (C.2) the normal density of Equation (C.1) we yield

$$2486 d_{\text{KL}}(Q, P) = \mathbb{E}_Q \left[ \log_e \left( \frac{1}{(2\pi)^{b/2} |\boldsymbol{\Sigma}_Q|^{1/2}} \right) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_Q)^\top \boldsymbol{\Sigma}_Q^{-1}(\mathbf{x} - \boldsymbol{\mu}_Q) \right. \\ 2487 \left. - \left[ \log_e \left( \frac{1}{(2\pi)^{b/2} |\boldsymbol{\Sigma}_P|^{1/2}} \right) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_P)^\top \boldsymbol{\Sigma}_P^{-1}(\mathbf{x} - \boldsymbol{\mu}_P) \right] \right] \\ 2488 = \mathbb{E}_Q \left[ -\frac{b}{2} \log_e(2\pi) - \frac{1}{2} \log_e(|\boldsymbol{\Sigma}_Q|) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_Q)^\top \boldsymbol{\Sigma}_Q^{-1}(\mathbf{x} - \boldsymbol{\mu}_Q) + \frac{b}{2} \log_e(2\pi) \right. \\ 2489 \left. + \frac{1}{2} \log_e(|\boldsymbol{\Sigma}_P|) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_P)^\top \boldsymbol{\Sigma}_P^{-1}(\mathbf{x} - \boldsymbol{\mu}_P) \right] \\ 2490 = \frac{1}{2} \mathbb{E}_Q [\log_e(|\boldsymbol{\Sigma}_P|) - \log_e(|\boldsymbol{\Sigma}_Q|) + (\mathbf{x} - \boldsymbol{\mu}_P)^\top \boldsymbol{\Sigma}_P^{-1}(\mathbf{x} - \boldsymbol{\mu}_P) \\ 2491 \left. - (\mathbf{x} - \boldsymbol{\mu}_Q)^\top \boldsymbol{\Sigma}_Q^{-1}(\mathbf{x} - \boldsymbol{\mu}_Q)]. \quad (\text{C.3}) \right.$$

2493 The expected value of a constant,  $\mathbb{E}[c]$ , is equal to  $c$  itself and, thus, Equation (C.3) becomes

$$2494 d_{\text{KL}}(Q, P) = \frac{1}{2} \log_e \left( \frac{|\boldsymbol{\Sigma}_P|}{|\boldsymbol{\Sigma}_Q|} \right) + \frac{1}{2} \mathbb{E}_Q [(\mathbf{x} - \boldsymbol{\mu}_P)^\top \boldsymbol{\Sigma}_P^{-1}(\mathbf{x} - \boldsymbol{\mu}_P) - (\mathbf{x} - \boldsymbol{\mu}_Q)^\top \boldsymbol{\Sigma}_Q^{-1}(\mathbf{x} - \boldsymbol{\mu}_Q)] \\ 2495 = \frac{1}{2} \log_e(|\boldsymbol{\Sigma}_Q^{-1} \boldsymbol{\Sigma}_P|) + \frac{1}{2} \mathbb{E}_Q [(\mathbf{x} - \boldsymbol{\mu}_P)^\top \boldsymbol{\Sigma}_P^{-1}(\mathbf{x} - \boldsymbol{\mu}_P)] \\ 2496 - \frac{1}{2} \mathbb{E}_Q [(\mathbf{x} - \boldsymbol{\mu}_Q)^\top \boldsymbol{\Sigma}_Q^{-1}(\mathbf{x} - \boldsymbol{\mu}_Q)]. \quad (\text{C.4})$$

The vector-matrix-vector product,  $(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})$ , produces a scalar whose value depends on the  
2498 entries of the vector  $\mathbf{x}$ . The covariance matrix,  $\Sigma$ , is constant and can be taken out of the expectation

$$\begin{aligned} 2499 \quad \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})] &= \mathbb{E}[\text{tr}((\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}))] \\ 2500 \quad &= \mathbb{E}[\text{tr}(\Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top)] \\ 2501 \quad &= \text{tr}(\mathbb{E}[\Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top]) \\ 2502 \quad &= \text{tr}(\Sigma^{-1}\mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top]), \\ 2503 \end{aligned} \tag{C.5}$$

2504 where the trace function

$$2505 \quad \text{tr}(\mathbf{A}) = \sum_{i=1}^b a_{ii} = a_{11} + a_{22} + \dots + a_{bb}, \\ 2506 \tag{C.6}$$

computes the sum of the elements on the main diagonal of the  $b \times b$  matrix,  $\Sigma^{-1}\mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top]$ .

2507 We can use Equation (C.5) to write Equation (C.4) as follows

$$\begin{aligned} 2508 \quad d_{\text{KL}}(Q, P) &= \frac{1}{2} \log_e(|\Sigma_Q^{-1}\Sigma_P|) + \frac{1}{2} \text{tr}(\Sigma_P^{-1}\mathbb{E}_Q[(\mathbf{x} - \boldsymbol{\mu}_P)(\mathbf{x} - \boldsymbol{\mu}_P)^\top]) \\ 2509 \quad &\quad - \frac{1}{2} \text{tr}(\Sigma_Q^{-1}\mathbb{E}_Q[(\mathbf{x} - \boldsymbol{\mu}_Q)(\mathbf{x} - \boldsymbol{\mu}_Q)^\top]), \\ 2510 \end{aligned} \tag{C.7}$$

2511 The expected value of the vector outer product,  $(\mathbf{x} - \boldsymbol{\mu}_Q)(\mathbf{x} - \boldsymbol{\mu}_Q)^\top$ , with respect to  $Q$  is simply equal to  
the covariance matrix,  $\Sigma_Q$ , of this distribution. Thus, we yield

$$\begin{aligned} 2512 \quad d_{\text{KL}}(Q, P) &= \frac{1}{2} \log_e(|\Sigma_Q^{-1}\Sigma_P|) - \frac{1}{2} \text{tr}(\Sigma_Q^{-1}\Sigma_Q) + \frac{1}{2} \text{tr}(\Sigma_P^{-1}\mathbb{E}_Q[(\mathbf{x} - \boldsymbol{\mu}_P)(\mathbf{x} - \boldsymbol{\mu}_P)^\top]) \\ 2513 \quad &= \frac{1}{2} \log_e(|\Sigma_Q^{-1}\Sigma_P|) - \frac{1}{2}b + \frac{1}{2} \text{tr}(\Sigma_P^{-1}\mathbb{E}_Q[\mathbf{x}\mathbf{x}^\top - \mathbf{x}\boldsymbol{\mu}_P^\top - \boldsymbol{\mu}_P\mathbf{x}^\top + \boldsymbol{\mu}_P\boldsymbol{\mu}_P^\top]) \\ 2514 \quad &= \frac{1}{2} \log_e(|\Sigma_Q^{-1}\Sigma_P|) - \frac{1}{2}b + \frac{1}{2} \text{tr}(\Sigma_P^{-1}\mathbb{E}_Q[\mathbf{x}\mathbf{x}^\top - 2\mathbf{x}\boldsymbol{\mu}_P^\top + \boldsymbol{\mu}_P\boldsymbol{\mu}_P^\top]), \\ 2515 \end{aligned} \tag{C.8}$$

2516 where the sum of the diagonal elements of the  $b \times b$  identity matrix,  $\mathbf{I}_b = \Sigma^{-1}\Sigma$ , equals the dimension,  
 $b \in \mathbb{N}_+$ , of the multivariate normal distribution and the  $b \times b$  matrix,  $-\mathbf{x}\boldsymbol{\mu}_P^\top - \boldsymbol{\mu}_P\mathbf{x}^\top$ , is conveniently  
written as  $-2\mathbf{x}\boldsymbol{\mu}_P^\top$ . This formulation for the sum of the two vector outer products holds only for the  
main diagonal elements of  $\mathbf{x}\boldsymbol{\mu}_P^\top - \boldsymbol{\mu}_P\mathbf{x}^\top$  on which the trace function operates.

2517

We can generate a large collection of  $N$  points,  $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$ , from  $Q \sim \mathcal{N}_b(\boldsymbol{\mu}_Q, \Sigma_Q)$ , and compute  
numerically the expected (mean) value of the  $b \times b$  matrix,  $\mathbf{x}\mathbf{x}^\top - 2\mathbf{x}\boldsymbol{\mu}_P^\top + \boldsymbol{\mu}_P\boldsymbol{\mu}_P^\top$ , between square brackets  
of Equation (C.8). With a little bit more effort, however, we can yield an analytic expression for the  
KL-divergence. The expected value of  $\mathbf{x}$  under  $Q$  is equal to the mean,  $\boldsymbol{\mu}_Q$ , of this distribution. From  
the general definition of the covariance matrix, we can derive a simple expression for the expected value  
2518 of the vector outer product,  $\mathbf{x}\mathbf{x}^\top$ , in Equation (C.8) as follows

$$\begin{aligned} 2519 \quad \Sigma &= \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] \\ 2520 \quad &= \mathbb{E}[\mathbf{x}\mathbf{x}^\top - \mathbf{x}\boldsymbol{\mu}^\top - \boldsymbol{\mu}\mathbf{x}^\top + \boldsymbol{\mu}\boldsymbol{\mu}^\top] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}[\mathbf{x}\mathbf{x}^\top] - \mathbb{E}[\mathbf{x}\boldsymbol{\mu}^\top] - \mathbb{E}[\boldsymbol{\mu}\mathbf{x}^\top] + \mathbb{E}[\boldsymbol{\mu}\boldsymbol{\mu}^\top] \\
&= \mathbb{E}[\mathbf{x}\mathbf{x}^\top] - \mathbb{E}[\mathbf{x}]\boldsymbol{\mu}^\top - \boldsymbol{\mu}\mathbb{E}[\mathbf{x}^\top] + \boldsymbol{\mu}\boldsymbol{\mu}^\top \\
&= \mathbb{E}[\mathbf{x}\mathbf{x}^\top] - \boldsymbol{\mu}\boldsymbol{\mu}^\top - \boldsymbol{\mu}\boldsymbol{\mu}^\top + \boldsymbol{\mu}\boldsymbol{\mu}^\top \\
&= \mathbb{E}[\mathbf{x}\mathbf{x}^\top] - \boldsymbol{\mu}\boldsymbol{\mu}^\top \\
&\implies \mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \Sigma + \boldsymbol{\mu}\boldsymbol{\mu}^\top. \tag{C.9}
\end{aligned}$$

2527 Thus, Equation (C.8) is equal to

$$\begin{aligned}
d_{\text{KL}}(Q, P) &= \frac{1}{2} \log_e(|\Sigma_Q^{-1}\Sigma_P|) - \frac{1}{2}b + \frac{1}{2} \text{tr}(\Sigma_P^{-1}(\Sigma_Q + \boldsymbol{\mu}_Q\boldsymbol{\mu}_Q^\top - 2\boldsymbol{\mu}_Q\boldsymbol{\mu}_P^\top + \boldsymbol{\mu}_P\boldsymbol{\mu}_P^\top)) \\
&= \frac{1}{2} [\log_e(|\Sigma_Q^{-1}\Sigma_P|) - b + \text{tr}(\Sigma_P^{-1}\Sigma_Q + \Sigma_P^{-1}\boldsymbol{\mu}_Q\boldsymbol{\mu}_Q^\top - 2\Sigma_P^{-1}\boldsymbol{\mu}_Q\boldsymbol{\mu}_P^\top + \Sigma_P^{-1}\boldsymbol{\mu}_P\boldsymbol{\mu}_P^\top)] \\
&= \frac{1}{2} [\log_e(|\Sigma_Q^{-1}\Sigma_P|) - b + \text{tr}(\Sigma_P^{-1}\Sigma_Q) + \text{tr}(\Sigma_P^{-1}\boldsymbol{\mu}_Q\boldsymbol{\mu}_Q^\top) - 2\text{tr}(\Sigma_P^{-1}\boldsymbol{\mu}_Q\boldsymbol{\mu}_P^\top) \\
&\quad + \text{tr}(\Sigma_P^{-1}\boldsymbol{\mu}_P\boldsymbol{\mu}_P^\top)]. \tag{C.10}
\end{aligned}$$

2533 As corollary of Equation (C.5), we yield

$$2534 \quad \text{tr}(\Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top) = (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}), \tag{C.11}$$

2536 and, consequently, Equation (C.10) may be written as follows

$$2537 \quad d_{\text{KL}}(Q, P) = \frac{1}{2} [\log_e(|\Sigma_Q^{-1}\Sigma_P|) - b + \text{tr}(\Sigma_P^{-1}\Sigma_Q) + \boldsymbol{\mu}_Q^\top \Sigma_P^{-1} \boldsymbol{\mu}_Q - 2\boldsymbol{\mu}_Q^\top \Sigma_P^{-1} \boldsymbol{\mu}_P + \boldsymbol{\mu}_P^\top \Sigma_P^{-1} \boldsymbol{\mu}_P]. \tag{C.12}$$

2539 The three vector-matrix-vector products in the above expression can be factorized to yield

$$2540 \quad d_{\text{KL}}(Q, P) = \frac{1}{2} [\log_e(|\Sigma_Q^{-1}\Sigma_P|) - b + \text{tr}(\Sigma_P^{-1}\Sigma_Q) + (\boldsymbol{\mu}_Q - \boldsymbol{\mu}_P)^\top \Sigma_P^{-1}(\boldsymbol{\mu}_Q - \boldsymbol{\mu}_P)]. \tag{C.13}$$

2542 This concludes our derivation of the KL-divergence of two multivariate normal distributions,  $Q \sim \mathcal{N}_b(\boldsymbol{\mu}_Q, \Sigma_Q)$  and  $P = \mathcal{N}_b(\boldsymbol{\mu}_P, \Sigma_P)$  in  $\mathbb{R}^b$  with  $b \in \mathbb{N}_+$ .

2544 In the special case of  $b = 1$  or two univariate normal distributions,  $Q \sim \mathcal{N}(\mu_Q, \sigma_Q^2)$  and  $P = \mathcal{N}(\mu_P, \sigma_P^2)$ ,  
2545 the above expression simplifies to

$$\begin{aligned}
d_{\text{KL}}(Q, P) &= \frac{1}{2} \left[ \log_e \left( \frac{\sigma_P^2}{\sigma_Q^2} \right) - 1 + \frac{\sigma_Q^2}{\sigma_P^2} + \frac{(\mu_Q - \mu_P)^2}{\sigma_P^2} \right] \\
&= \log_e \left( \frac{\sigma_P}{\sigma_Q} \right) + \frac{\sigma_Q^2 + (\mu_Q - \mu_P)^2 - \sigma_P^2}{2\sigma_P^2}. \tag{C.14}
\end{aligned}$$

2549 The use of the natural logarithm in our derivation of Equations (C.13) and (C.14) affixes the unit of  
2550 nats to  $d_{\text{KL}}(Q, P)$ . To change units of the relative entropy, one should just divide  $d_{\text{KL}}(Q, P)$  by  $\log_e(z)$ .  
2551 Then for  $z = 2$  we yield  $d_{\text{KL}}(Q, P)$  in bits. Furthermore, we obtain the reverse KL-divergence  $d_{\text{KL}}(P, Q)$   
2552 by swapping arguments  $Q$  and  $P$  in the respective Equations.

## 2553 Appendix D: KL divergence of normal and uniform distributions

2554 We derive an analytic expression for the relative entropy  $d_{\text{KL}}(Q, P)$

$$2555 \quad d_{\text{KL}}(Q, P) = \mathbb{H}(Q, P) - \mathbb{H}(Q), \quad (\text{D.1})$$

2556 of a normal *true distribution*  $Q = \mathcal{N}(\mu_Q, \sigma_Q^2)$  and uniform *forecast distribution*  $P = \mathcal{U}(a_P, b_P)$  on a  
2557 bounded sample space  $x \in [a_P, b_P]$  and  $b_P > a_P$ . The cross-entropy of  $Q$  and  $P$  is equal to

$$\begin{aligned} 2558 \quad \mathbb{H}(Q, P) &= - \int_{a_P}^{b_P} Q(x) \log_b(P(x)) dx \\ 2559 \quad &= - \int_{a_P}^{b_P} Q(x) \log_e\left(\frac{1}{b_P - a_P}\right) dx \\ 2560 \quad &= \log_e(b_P - a_P) \int_{a_P}^{b_P} Q(x) dx \\ 2561 \quad &= \log_e(b_P - a_P), \end{aligned} \quad (\text{D.2})$$

2563 in units of nats. Note that the cross-entropy will attain an infinite value on the extended real line,  $x \in \overline{\mathbb{R}}$   
2564 as  $P(x) \rightarrow 0$ . The entropy of the normal *true distribution*  $Q$  may be computed as follows

$$\begin{aligned} 2565 \quad \mathbb{H}(Q) &= - \int_{a_P}^{b_P} Q(x) \log_b(Q(x)) dx \\ 2566 \quad &= - \int_{a_P}^{b_P} Q(x) \log_e\left[\frac{1}{\sigma_Q \sqrt{2\pi}} \exp\left(-\frac{1}{2} \frac{(x - \mu_Q)^2}{\sigma_Q^2}\right)\right] dx \\ 2567 \quad &= - \int_{a_P}^{b_P} Q(x) \left(\log_e\left[\frac{1}{\sigma_Q \sqrt{2\pi}}\right] - \frac{1}{2} \frac{(x - \mu_Q)^2}{\sigma_Q^2}\right) dx \\ 2568 \quad &= \log_e(2\pi\sigma_Q^2)^{1/2} \int_{a_P}^{b_P} Q(x) dx + \frac{1}{2\sigma_Q^2} \int_{a_P}^{b_P} (x - \mu_Q)^2 Q(x) dx \\ 2569 \quad &= \frac{1}{2} \log_e(2\pi\sigma_Q^2) \times 1 + \frac{1}{2\sigma_Q^2} \times \sigma_Q^2 \\ 2570 \quad &= \frac{1}{2} \log_e(2e\pi\sigma_Q^2). \end{aligned} \quad (\text{D.3})$$

2572 The relative entropy is now equal to

$$\begin{aligned} 2573 \quad d_{\text{KL}}(Q, P) &= \mathbb{H}(Q, P) - \mathbb{H}(Q) \\ 2574 \quad &= \log_e(b_P - a_P) - \frac{1}{2} \log_e(2e\pi\sigma_Q^2), \end{aligned} \quad (\text{D.4})$$

2576 in units of nats. If the interval of  $P(x)$  exceeds by far that of  $Q(x)$  and, thus,  $b_P - a_P \gg \sigma_Q$  then  
2577  $d_{\text{KL}}(Q, P)$  may attain an infinite value. For any scalar  $\beta \in \mathbb{R}_+$  and  $\sigma_Q = (b_P - a_P)/\nu$  the relative entropy  
2578 simplifies to  $d_{\text{KL}}(Q, P) = -\frac{1}{2} \log_e(2e\pi) + \log_e(\nu)$ .

2579 To resolve problems with the uniform distribution of  $P$  on an unbounded interval we could specify  
 2580  $P(x) \propto 1$  instead. Then the cross-entropy  $\mathbb{H}(Q, P) = 0$  and the relative entropy  $d_{\text{KL}}(Q, P)$  reduces to  
 2581 the so-called differential entropy  $\frac{1}{2} \log_e(2e\pi\sigma_Q^2)$  of the normal distribution  $Q$ .

2582 We can follow a similar derivation for the reverse KL-divergence,  $d_{\text{KL}}(P, Q)$ . The cross-entropy of  $P$  and  
 2583  $Q$  in units of nats is equal to

$$\begin{aligned}
 2584 \quad \mathbb{H}(P, Q) &= - \int_{a_P}^{b_P} P(x) \log_b(Q(x)) dx \\
 2585 &= - \int_{a_P}^{b_P} P(x) \log_b \left[ \frac{1}{\sigma_Q \sqrt{2\pi}} \exp \left( -\frac{1}{2} \frac{(x - \mu_Q)^2}{\sigma_Q^2} \right) \right] dx \\
 2586 &= - \int_{a_P}^{b_P} P(x) \left( \log_e \left[ \frac{1}{\sigma_Q \sqrt{2\pi}} \right] - \frac{1}{2} \frac{(x - \mu_Q)^2}{\sigma_Q^2} \right) dx \\
 2587 &= \log_e(2\pi\sigma_Q^2)^{1/2} \int_{a_P}^{b_P} P(x) dx + \frac{1}{2\sigma_Q^2} \int_{a_P}^{b_P} (x - \mu_Q)^2 P(x) dx \\
 2588 &= \frac{1}{2} \log_e(2\pi\sigma_Q^2) + \frac{1}{2(b_P - a_P)\sigma_Q^2} \left| -\frac{1}{3}(\mu_Q - x)^3 \right|_{a_P}^{b_P} \\
 2589 &= \frac{1}{2} \log_e(2\pi\sigma_Q^2) + \frac{(\mu_Q - a_P)^3 - (\mu_Q - b_P)^3}{6(b_P - a_P)\sigma_Q^2}. \\
 2590
 \end{aligned} \tag{D.5}$$

2591 The entropy of  $P$  is equal to Equation (D.2) to yield

$$\begin{aligned}
 2592 \quad \mathbb{H}(P) &= - \int_{a_P}^{b_P} P(x) \log_b(P(x)) dx \\
 2593 &= - \int_{a_P}^{b_P} P(x) \log_e \left( \frac{1}{b_P - a_P} \right) dx \\
 2594 &= \log_e(b_P - a_P) \int_{a_P}^{b_P} P(x) dx \\
 2595 &= \log_e(b_P - a_P), \\
 2596
 \end{aligned} \tag{D.6}$$

2597 in units of nats. Now we yield the following expression for the relative entropy  $d_{\text{KL}}(P, Q)$  in nats

$$\begin{aligned}
 2598 \quad d_{\text{KL}}(P, Q) &= \mathbb{H}(P, Q) - \mathbb{H}(P) \\
 2599 &= \frac{1}{2} \log_e(2\pi\sigma_Q^2) + \frac{(\mu_Q - a_P)^3 - (\mu_Q - b_P)^3}{6(b_P - a_P)\sigma_Q^2} - \log_e(b_P - a_P). \\
 2600
 \end{aligned} \tag{D.7}$$

2601 This confirms again that the relative entropy is not symmetric in  $Q$  and  $P$ .

## 2602 Appendix E: Rainfall data

2603 Table E.1 is taken from *Hughes and Topp* (2015) and summarizes a data set of  $n = 346$  forecasts of  
 2604 24-hour precipitation probability made by the Finnish Meteorological Institute during 2003 for the city  
 2605 of Tampere in south-central Finland. The left block presents the original data, and the right block lists  
 the data used in our case study.

Table E.1: Rainfall data from Table 1 of *Hughes and Topp* (2015) for the city of Tampere, Finland.

a: Original data						b: Adapted data		
$k$	$p_k$	$n_k$	$o_k$	$\bar{o}_k$	$n_k/n$	$k$	$p_k$	$q_k$
1	0.05	46	1	0.0217	0.1329	1	0.1727	0.1359
2	0.1	55	1	0.0182	0.1590	2	0.1636	0.1364
3	0.2	59	5	0.0847	0.1705	3	0.1455	0.1272
4	0.3	41	5	0.1220	0.1185	4	0.1273	0.1220
5	0.4	19	4	0.2105	0.0549	5	0.1091	0.1097
6	0.5	22	8	0.3636	0.0636	6	0.0909	0.0884
7	0.6	22	6	0.2727	0.0636	7	0.0727	0.1011
8	0.7	34	16	0.4706	0.0983	8	0.0545	0.0736
9	0.8	24	16	0.6667	0.0694	9	0.0364	0.0463
10	0.9	11	8	0.7273	0.0318	10	0.0182	0.0379
11	0.95	13	11	0.8462	0.0376	11	0.0091	0.0214
$\sum$		346	81		1.0000	$\sum$		1.000

2606

2607 Forecast probabilities of rainfall  $p_k$ ;  $k = (1, \dots, m)$  were issued using  $m = 11$  categories. The variable  $n_k$   
 2608 lists the number of days for which the Finnish Meteorological Institute quoted  $p_k$ . Then,  $o_k$ , signifies the  
 2609 number of days on which measured rainfall depths for the city of Tampere exceeded  $\geq 0.2$  mm, otherwise a  
 2610 no-rain day was recorded. Next, the ratio  $\bar{o}_k = o_k/n_k$  equals the *true* rainfall probability for each forecast  
 2611 category. Finally,  $n_k/n$  corresponds to the relative frequency of each forecast category. We refer readers  
 2612 to *Hughes and Topp* (2015) for a more detailed description of the data set and its use in a diagrammatic  
 2613 interpretation of the Brier scoring rule and related score divergence. The raw precipitation data can be  
 2614 found at [https://www.cawcr.gov.au/projects/verification/POP3/POP\\_3cat\\_2003.txt](https://www.cawcr.gov.au/projects/verification/POP3/POP_3cat_2003.txt).

2615 The right block tabulates the data that was used in our case study. The forecast probabilities of the  
 2616 individual categories are normalized to sum to unity. This defines a  $m$ -vector  $\mathbf{p} = (p_1, \dots, p_m)^\top$  of rainfall  
 2617 probability forecasts. We apply a similar normalization to the  $\bar{o}_k$ 's to yield the vector  $\mathbf{q} = (q_1, \dots, q_m)^\top$   
 2618 of *true* rainfall probabilities.

## 2619 Appendix F: Quantile form of Continuous Ranked Probability Score

2620 The CRPS of a distribution forecast  $P$  and verifying observation  $\omega \in \Omega$  is equal to the integral of the  
 2621 quantile score function

$$2622 \quad S_{\text{CRPS}}(P, \omega) = \int_0^1 S_{\text{QNT}}^\tau(P, \omega) d\tau, \quad (F.1)$$

2623

2624 and, thus, we yield

$$2625 \quad S_{\text{CRPS}}(P, \omega) = -2 \int_0^1 (\mathbb{1}\{\omega < y_\tau\} - \tau)(y_\tau - \omega) d\tau. \quad (F.2)$$

2626

2627 The above expression may be rearranged and written as follows

$$\begin{aligned} 2628 \quad S_{\text{CRPS}}(P, \omega) &= 2 \int_0^1 \tau(y_\tau - \omega) d\tau - 2 \int_0^1 \mathbb{1}\{\omega < y_\tau\}(y_\tau - \omega) d\tau \\ 2629 &= 2 \int_0^1 \tau(F_P^{-1}(\tau) - \omega) d\tau - 2 \int_0^1 \mathbb{1}\{F_P(\omega) < \tau\}(F_P^{-1}(\tau) - \omega) d\tau \\ 2630 &= 2 \int_0^1 \tau F_P^{-1}(\tau) d\tau - 2\omega \int_0^1 \tau d\tau - 2 \int_{F_P(\omega)}^1 (F_P^{-1}(\tau) - \omega) d\tau \\ 2631 &= 2 \int_0^1 \tau F_P^{-1}(\tau) d\tau - 2\omega \int_0^1 \tau d\tau - 2 \int_{F_P(\omega)}^1 F_P^{-1}(\tau) d\tau + 2\omega \int_{F_P(\omega)}^1 d\tau \\ 2632 &= 2 \int_0^1 \tau F_P^{-1}(\tau) d\tau - 2\omega \left[ \frac{1}{2}\tau^2 \right]_0^1 - 2 \int_{F_P(\omega)}^1 F_P^{-1}(\tau) d\tau + 2\omega \left[ \tau \right]_{F_P(\omega)}^1 \\ 2633 &= 2 \int_0^1 \tau F_P^{-1}(\tau) d\tau - \omega - 2 \int_{F_P(\omega)}^1 F_P^{-1}(\tau) d\tau + 2\omega - 2\omega F_P(\omega) \\ 2634 &= 2 \int_0^1 \tau F_P^{-1}(\tau) d\tau + \omega - 2 \int_{F_P(\omega)}^1 F_P^{-1}(\tau) d\tau - 2\omega F_P(\omega), \end{aligned} \quad (F.3)$$

2635

2636 and finally, we yield

$$2637 \quad S_{\text{CRPS}}(P, \omega) = \omega(1 - 2F_P(\omega)) + 2 \int_0^1 \tau F_P^{-1}(\tau) d\tau - 2 \int_{F_P(\omega)}^1 F_P^{-1}(\tau) d\tau. \quad (F.4)$$

2638

2639 This concludes the derivation.

## 2640 Appendix G: Continuous Ranked Probability Score for $\mathcal{N}(\mu, \sigma^2)$

2641 The continuous ranked probability score (CRPS) is given by

$$2642 S_{\text{CRPS}}(P, \omega) = - \int_{-\infty}^{\infty} (F_P(z) - \mathbb{1}\{z \geq \omega\})^2 dz, \quad (\text{G.1})$$

2643 where  $F_P(z)$  denotes the cumulative distribution function of  $P$  and the indicator function,  $\mathbb{1}\{a\}$ , returns  
 2644 1 if  $a$  is true and zero otherwise. If we make the convenient assumption that the probability measure  
 2645 is univariate normal,  $P = \mathcal{N}(\mu, \sigma^2)$ , then the cumulative distribution function of  $P$  has a closed-form  
 2646 expression

$$2647 F_P(x, \mu, \sigma^2) = \frac{1}{2} \left[ 1 + \operatorname{erf} \left( \frac{x - \mu}{\sigma \sqrt{2}} \right) \right], \quad (\text{G.2})$$

2648 where  $\operatorname{erf}(x)$  is the error function for element  $x$

$$2649 \operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x \exp(-t^2) dt, \quad (\text{G.3})$$

2650 and the CRPS becomes

$$2651 S_{\text{CRPS}}(\mathcal{N}(\mu, \sigma^2), \omega) = - \int_{-\infty}^{\omega} \left( \frac{1}{2} \left[ 1 + \operatorname{erf} \left( \frac{z - \mu}{\sigma \sqrt{2}} \right) \right] \right)^2 dz - \int_{\omega}^{\infty} \left( \frac{1}{2} \left[ 1 + \operatorname{erf} \left( \frac{z - \mu}{\sigma \sqrt{2}} \right) \right] - 1 \right)^2 dz. \quad (\text{G.4})$$

2653 We use the symbolic toolbox in MATLAB to yield closed-form expressions for the two definite integrals  
 2654 in the above expression. We first write out the left integral

$$\begin{aligned} 2655 - \int_{-\infty}^{\omega} \left( \frac{1}{2} \left[ 1 + \operatorname{erf} \left( \frac{z - \mu}{\sigma \sqrt{2}} \right) \right] \right)^2 dz &= - \left| \frac{z}{4} + \frac{(\mu - z)}{2} \operatorname{erf} \left( \frac{\mu - z}{\sigma \sqrt{2}} \right) - \frac{(\mu - z)}{4} \operatorname{erf} \left( \frac{\mu - z}{\sigma \sqrt{2}} \right)^2 \right. \\ 2656 &\quad \left. + \frac{\sigma}{2\sqrt{\pi}} \operatorname{erf} \left( \frac{\mu - z}{\sigma} \right) + \frac{\sigma}{\sqrt{2\pi}} \exp \left( -\frac{(\mu - z)^2}{2\sigma^2} \right) \right. \\ 2657 &\quad \left. - \frac{\sigma}{\sqrt{2\pi}} \exp \left( -\frac{(\mu - z)^2}{2\sigma^2} \right) \operatorname{erf} \left( \frac{\mu - z}{\sigma \sqrt{2}} \right) \right|_{-\infty}^{\omega} \\ 2658 &= \left| -\frac{z}{4} - \frac{(\mu - z)}{2} f(z, \mu, \sigma) + \frac{(\mu - z)}{4} f(z, \mu, \sigma)^2 - \frac{\sigma}{2\sqrt{\pi}} g(z, \mu, \sigma) \right. \\ 2659 &\quad \left. - \frac{\sigma}{\sqrt{2\pi}} h(z, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} f(z, \mu, \sigma) h(z, \mu, \sigma) \right|_{-\infty}^{\omega}, \end{aligned} \quad (\text{G.5})$$

2661 where  $f(z, \mu, \sigma) = \operatorname{erf}\left(\frac{1}{2}\sqrt{2}(\mu - z)/\sigma\right)$ ,  $g(z, \mu, \sigma) = \operatorname{erf}\left((\mu - z)/\sigma\right)$  and  $h(z, \mu, \sigma) = \exp\left(-\frac{1}{2}(\mu - z)^2/\sigma^2\right)$ .

2662 We follow a similar recipe for the right integral of Equation (G.4) to yield

$$\begin{aligned}
 2663 \quad & - \int_{-\infty}^{\omega} \left( \frac{1}{2} \left[ 1 + \operatorname{erf}\left(\frac{\omega - \mu}{\sigma\sqrt{2}}\right) \right] - 1 \right)^2 dz = - \left| \frac{z}{4} - \frac{(\mu - z)}{2} \operatorname{erf}\left(\frac{\mu - z}{\sigma\sqrt{2}}\right) - \frac{(\mu - z)}{4} \operatorname{erf}\left(\frac{\mu - z}{\sigma\sqrt{2}}\right)^2 \right. \\
 2664 \quad & \quad \left. + \frac{\sigma}{2\sqrt{\pi}} \operatorname{erf}\left(\frac{\mu - z}{\sigma}\right) - \frac{\sigma}{\sqrt{2\pi}} \exp\left(-\frac{(\mu - z)^2}{2\sigma^2}\right) \right. \\
 2665 \quad & \quad \left. - \frac{\sigma}{\sqrt{2\pi}} \exp\left(-\frac{(\mu - z)^2}{2\sigma^2}\right) \operatorname{erf}\left(\frac{\mu - z}{\sigma\sqrt{2}}\right) \right|_{-\infty}^{\omega} \\
 2666 \quad & = \left| -\frac{z}{4} + \frac{(\mu - z)}{2} f(z, \mu, \sigma) + \frac{(\mu - z)}{4} f(z, \mu, \sigma)^2 - \frac{\sigma}{2\sqrt{\pi}} g(z, \mu, \sigma) \right. \\
 2667 \quad & \quad \left. + \frac{\sigma}{\sqrt{2\pi}} h(z, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} f(z, \mu, \sigma) h(z, \mu, \sigma) \right|_{-\infty}^{\omega}. \tag{G.6}
 \end{aligned}$$

2669 Before admitting the integral bounds, we first perform limit analysis of the constituent functions

$$\begin{aligned}
 2670 \quad & \lim_{z \rightarrow -\infty} f(z, \mu, \sigma) = \lim_{z \rightarrow -\infty} \operatorname{erf}\left(\frac{\mu - z}{\sigma\sqrt{2}}\right) = 1 \\
 2671 \quad & \lim_{z \rightarrow \infty} f(z, \mu, \sigma) = \lim_{z \rightarrow \infty} \operatorname{erf}\left(\frac{\mu - z}{\sigma\sqrt{2}}\right) = -1 \\
 2672 \quad & \lim_{z \rightarrow -\infty} g(z, \mu, \sigma) = \lim_{z \rightarrow -\infty} \operatorname{erf}\left(\frac{\mu - z}{\sigma}\right) = 1 \\
 2673 \quad & \lim_{z \rightarrow \infty} g(z, \mu, \sigma) = \lim_{z \rightarrow \infty} \operatorname{erf}\left(\frac{\mu - z}{\sigma}\right) = -1 \\
 2674 \quad & \lim_{z \rightarrow -\infty} h(z, \mu, \sigma) = \lim_{z \rightarrow -\infty} \exp\left(-\frac{(\mu - z)^2}{2\sigma^2}\right) = 0 \\
 2675 \quad & \lim_{z \rightarrow \infty} h(z, \mu, \sigma) = \lim_{z \rightarrow \infty} \exp\left(-\frac{(\mu - z)^2}{2\sigma^2}\right) = 0. \tag{G.7}
 \end{aligned}$$

2677 The left integral is now equal to

$$\begin{aligned}
 2678 \quad & - \int_{-\infty}^{\omega} \left( \frac{1}{2} \left[ 1 + \operatorname{erf}\left(\frac{z - \mu}{\sigma\sqrt{2}}\right) \right] \right)^2 dz = \left| -\frac{z}{4} - \frac{(\mu - z)}{2} f(z, \mu, \sigma) + \frac{(\mu - z)}{4} f(z, \mu, \sigma)^2 - \frac{\sigma}{2\sqrt{\pi}} g(z, \mu, \sigma) \right. \\
 2679 \quad & \quad \left. - \frac{\sigma}{\sqrt{2\pi}} h(z, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} f(z, \mu, \sigma) h(z, \mu, \sigma) \right|_{-\infty}^{\omega} \\
 2680 \quad & = \left( -\frac{\omega}{4} - \frac{(\mu - \omega)}{2} f(\omega, \mu, \sigma) + \frac{(\mu - \omega)}{4} f(\omega, \mu, \sigma)^2 - \frac{\sigma}{2\sqrt{\pi}} g(\omega, \mu, \sigma) \right. \\
 2681 \quad & \quad \left. - \frac{\sigma}{\sqrt{2\pi}} h(\omega, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} f(\omega, \mu, \sigma) h(\omega, \mu, \sigma) \right) \\
 2682 \quad & \quad - \lim_{z \rightarrow -\infty} \left( -\frac{z}{4} - \frac{(\mu - z)}{2} f(z, \mu, \sigma) + \frac{(\mu - z)}{4} f(z, \mu, \sigma)^2 \right. \\
 2683 \quad & \quad \left. - \frac{\sigma}{2\sqrt{\pi}} g(z, \mu, \sigma) - \frac{\sigma}{\sqrt{2\pi}} h(z, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} f(z, \mu, \sigma) h(z, \mu, \sigma) \right) \\
 2684 \quad & = \left( -\frac{\omega}{4} - \frac{(\mu - \omega)}{2} f(\omega, \mu, \sigma) + \frac{(\mu - \omega)}{4} f(\omega, \mu, \sigma)^2 - \frac{\sigma}{2\sqrt{\pi}} g(\omega, \mu, \sigma) \right)
 \end{aligned}$$

$$\begin{aligned}
& - \frac{\sigma}{\sqrt{2\pi}} h(\omega, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} f(\omega, \mu, \sigma) h(\omega, \mu, \sigma) \Big) \\
& - \lim_{z \rightarrow -\infty} \left( -\frac{z}{4} - \frac{(\mu - z)}{2} \times 1 + \frac{(\mu - z)}{4} \times 1^2 - \frac{\sigma}{2\sqrt{\pi}} \times 1 \right. \\
& \quad \left. - \frac{\sigma}{\sqrt{2\pi}} \times 0 + \frac{\sigma}{\sqrt{2\pi}} \times 0 \times 1 \right) \\
= & \left( -\frac{\omega}{4} - \frac{(\mu - \omega)}{2} f(\omega, \mu, \sigma) + \frac{(\mu - \omega)}{4} f(\omega, \mu, \sigma)^2 - \frac{\sigma}{2\sqrt{\pi}} g(\omega, \mu, \sigma) \right. \\
& \quad \left. - \frac{\sigma}{\sqrt{2\pi}} h(\omega, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} f(\omega, \mu, \sigma) h(\omega, \mu, \sigma) \right) \\
& - \lim_{z \rightarrow -\infty} \left( -\frac{z}{4} + \frac{z}{2} - \frac{\mu}{2} - \frac{z}{4} + \frac{\mu}{4} - \frac{\sigma}{2\sqrt{\pi}} \right) \\
= & -\frac{\omega}{4} - \frac{(\mu - \omega)}{2} f(\omega, \mu, \sigma) + \frac{(\mu - \omega)}{4} f(\omega, \mu, \sigma)^2 - \frac{\sigma}{2\sqrt{\pi}} g(\omega, \mu, \sigma) \\
& - \frac{\sigma}{\sqrt{2\pi}} h(\omega, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} f(\omega, \mu, \sigma) h(\omega, \mu, \sigma) + \frac{\mu}{4} + \frac{\sigma}{2\sqrt{\pi}}, \quad (\text{G.8})
\end{aligned}$$

and the right integral becomes

$$\begin{aligned}
& - \int_{\omega}^{\infty} \left( \frac{1}{2} \left[ 1 + \operatorname{erf} \left( \frac{z - \mu}{\sigma\sqrt{2}} \right) \right] - 1 \right)^2 dz = \left| -\frac{z}{4} + \frac{(\mu - z)}{2} f(z, \mu, \sigma) + \frac{(\mu - z)}{4} f(z, \mu, \sigma)^2 - \frac{\sigma}{2\sqrt{\pi}} g(z, \mu, \sigma) \right. \\
& \quad \left. + \frac{\sigma}{\sqrt{2\pi}} h(z, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} f(z, \mu, \sigma) h(z, \mu, \sigma) \right|_{\omega}^{\infty} \\
= & \lim_{z \rightarrow \infty} \left( -\frac{z}{4} + \frac{(\mu - z)}{2} f(z, \mu, \sigma) + \frac{(\mu - z)}{4} f(z, \mu, \sigma)^2 \right. \\
& \quad \left. - \frac{\sigma}{2\sqrt{\pi}} g(z, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} h(z, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} f(z, \mu, \sigma) h(z, \mu, \sigma) \right) \\
& - \left( -\frac{\omega}{4} + \frac{(\mu - \omega)}{2} f(\omega, \mu, \sigma) + \frac{(\mu - \omega)}{4} f(\omega, \mu, \sigma)^2 \right. \\
& \quad \left. - \frac{\sigma}{2\sqrt{\pi}} g(\omega, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} h(\omega, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} f(\omega, \mu, \sigma) h(\omega, \mu, \sigma) \right) \\
= & \lim_{z \rightarrow \infty} \left( -\frac{z}{4} + \frac{(\mu - z)}{2} \times (-1) + \frac{(\mu - z)}{4} \times (-1)^2 - \frac{\sigma}{2\sqrt{\pi}} \times (-1) \right. \\
& \quad \left. + \frac{\sigma}{\sqrt{2\pi}} \times 0 + \frac{\sigma}{\sqrt{2\pi}} \times 0 \times (-1) \right) - \left( -\frac{\omega}{4} + \frac{(\mu - \omega)}{2} f(\omega, \mu, \sigma) \right. \\
& \quad \left. + \frac{(\mu - \omega)}{4} f(\omega, \mu, \sigma)^2 - \frac{\sigma}{2\sqrt{\pi}} g(\omega, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} h(\omega, \mu, \sigma) \right. \\
& \quad \left. + \frac{\sigma}{\sqrt{2\pi}} f(\omega, \mu, \sigma) h(\omega, \mu, \sigma) \right) \\
= & \lim_{z \rightarrow \infty} \left( -\frac{z}{4} + \frac{z}{2} - \frac{\mu}{2} - \frac{z}{4} + \frac{\mu}{4} + \frac{\sigma}{2\sqrt{\pi}} \right) \\
& - \left( -\frac{\omega}{4} + \frac{(\mu - \omega)}{2} f(\omega, \mu, \sigma) + \frac{(\mu - \omega)}{4} f(\omega, \mu, \sigma)^2 \right. \\
& \quad \left. - \frac{\sigma}{2\sqrt{\pi}} g(\omega, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} h(\omega, \mu, \sigma) + \frac{\sigma}{\sqrt{2\pi}} f(\omega, \mu, \sigma) h(\omega, \mu, \sigma) \right)
\end{aligned}$$

2708

$$= -\frac{\mu}{4} + \frac{\sigma}{2\sqrt{\pi}} + \frac{\omega}{4} - \frac{(\mu - \omega)}{2} f(\omega, \mu, \sigma) - \frac{(\mu - \omega)}{4} f(\omega, \mu, \sigma)^2$$

2709

$$+ \frac{\sigma}{2\sqrt{\pi}} g(\omega, \mu, \sigma) - \frac{\sigma}{\sqrt{2\pi}} h(\omega, \mu, \sigma) - \frac{\sigma}{\sqrt{2\pi}} f(\omega, \mu, \sigma) h(\omega, \mu, \sigma).$$

2710 (G.9)

2711 We can now add up Equations (G.8) and (G.9) to yield the CRPS of  $P = \mathcal{N}(\mu, \sigma^2)$  in Equation (G.4)

2712  $S_{\text{CRPS}}(\mathcal{N}(\mu, \sigma^2), \omega) = \left( \cancel{-\frac{\omega}{4}} - \frac{(\mu - \omega)}{2} f(\omega, \mu, \sigma) + \cancel{\frac{(\mu - \omega)}{4} f(\omega, \mu, \sigma)^2} - \cancel{\frac{\sigma}{2\sqrt{\pi}} g(\omega, \mu, \sigma)}$

2713  $- \frac{\sigma}{\sqrt{2\pi}} h(\omega, \mu, \sigma) + \cancel{\frac{\sigma}{\sqrt{2\pi}} f(\omega, \mu, \sigma) h(\omega, \mu, \sigma)} + \cancel{\frac{\mu}{4}} + \frac{\sigma}{2\sqrt{\pi}} \right) +$

2714  $\left( \cancel{-\frac{\mu}{4}} + \frac{\sigma}{2\sqrt{\pi}} \cancel{+\frac{\omega}{4}} - \frac{(\mu - \omega)}{2} f(\omega, \mu, \sigma) - \cancel{\frac{(\mu - \omega)}{4} f(\omega, \mu, \sigma)^2}$

2715  $+ \cancel{\frac{\sigma}{2\sqrt{\pi}} g(\omega, \mu, \sigma)} - \cancel{\frac{\sigma}{\sqrt{2\pi}} h(\omega, \mu, \sigma)} - \cancel{\frac{\sigma}{\sqrt{2\pi}} f(\omega, \mu, \sigma) h(\omega, \mu, \sigma)} \right)$

2716  $= -\frac{2(\mu - \omega)}{2} f(\omega, \mu, \sigma) - \frac{2\sigma}{\sqrt{2\pi}} h(\omega, \mu, \sigma) + \frac{2\sigma}{2\sqrt{\pi}}$

2717  $= -(\mu - \omega) \text{erf}\left(\frac{\mu - \omega}{\sigma\sqrt{2}}\right) - \frac{\sigma\sqrt{2}}{\sqrt{\pi}} \exp\left(-\frac{(\mu - \omega)^2}{2\sigma^2}\right) + \frac{\sigma}{\sqrt{\pi}}.$

2718 (G.10)

2719 We can manipulate this expression into a function of the normal PDF,  $f_{\mathcal{N}}(x, \mu, \sigma^2)$ , and normal CDF,

2720  $F_{\mathcal{N}}(x, \mu, \sigma^2)$ , in Equation (G.2) as follows

2721  $S_{\text{CRPS}}(\mathcal{N}(\mu, \sigma^2), \omega) = -(\omega - \mu) \text{erf}\left(\frac{\omega - \mu}{\sigma\sqrt{2}}\right) - 2\sigma^2 \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(\mu - \omega)^2}{2\sigma^2}\right) + \frac{\sigma}{\sqrt{\pi}}$

2722  $= -(\omega - \mu) \left[ 1 + \text{erf}\left(\frac{\omega - \mu}{\sigma\sqrt{2}}\right) \right] + (\omega - \mu) - 2\sigma^2 f_{\mathcal{N}}(\omega, \mu, \sigma^2) + \frac{\sigma}{\sqrt{\pi}}$

2723  $= -2(\omega - \mu) \frac{1}{2} \left[ 1 + \text{erf}\left(\frac{\omega - \mu}{\sigma\sqrt{2}}\right) \right] + (\omega - \mu) - 2\sigma^2 f_{\mathcal{N}}(\omega, \mu, \sigma^2) + \frac{\sigma}{\sqrt{\pi}}$

2724  $= -2(\omega - \mu) F_{\mathcal{N}}(\omega, \mu, \sigma^2) + (\omega - \mu) - 2\sigma^2 f_{\mathcal{N}}(\omega, \mu, \sigma^2) + \frac{\sigma}{\sqrt{\pi}},$

2725 (G.11)

2726 which may be rearranged and simplified to

2727  $S_{\text{CRPS}}(\mathcal{N}(\mu, \sigma^2), \omega) = \frac{\sigma}{\sqrt{\pi}} - 2\sigma^2 f_{\mathcal{N}}(\omega, \mu, \sigma^2) - (\omega - \mu)(2F_{\mathcal{N}}(\omega, \mu, \sigma^2) - 1).$

2728 (G.12)

2729 Note that the quantile form of the CRPS in Equation (61) would lead to an equivalent solution as

2730 above but in fewer steps. This concludes the derivation of the CRPS for a normal distribution forecast

2731  $P = \mathcal{N}(\mu, \sigma^2)$  and verifying observation  $\omega \in \Omega$ .

## 2732 Appendix H: Description of Hydrologic models

2733 In this Appendix we present a brief overview of the HYMOD, Hmodel and SAC-SMA conceptual  
 2734 watershed models. The models are coded in MATLAB and C++ and use a mass-conservative second-  
 2735 order integration method with adaptive time step. This guarantees a robust and accurate numerical  
 2736 solution of the simulated fluxes and state variables. Next, we discuss each of the models separately.

### 2737 H.1 HYdrologic MODEL

2738 The HYdrologic MODel (HYMOD) originates from the PhD thesis of *Boyle* (2001) and describes the  
 2739 rainfall-discharge relationship using five fictitious control volumes. These reservoirs simulate processes  
 such as evaporation, percolation, river inflow and baseflow (see Figure H.1). Table H.1 presents the five

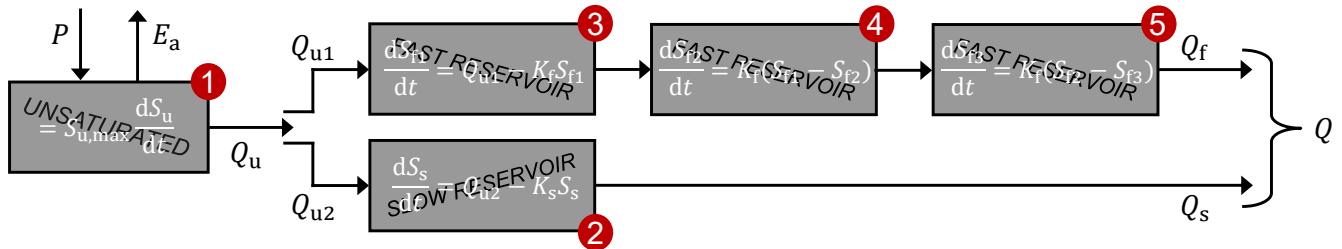


Figure H.1: Schematic illustration of the HYdrologic MODel of *Boyle* (2001). Grey boxes, labeled in red, correspond to fictitious control volumes of the watershed which govern the rainfall-runoff transformation. The state variables,  $S_u$ ,  $S_s$ ,  $S_{f1}$ ,  $S_{f2}$  and  $S_{f3}$ , correspond to the water storage in each compartment. Arrows portray the fluxes into and out of the compartments, including precipitation,  $P$ , evaporation,  $E_a$ , precipitation converted into flow,  $Q_u$ , fast flow,  $Q_f$ , and baseflow,  $Q_s$ . The fluxes are computed as follows,  $Q_u = P(1 - (1 - \bar{S}_u)^b)$ ,  $E_a = E_p \bar{S}_u (1 + c) / (\bar{S}_u + c)$ ,  $Q_{u1} = a Q_u$ ,  $Q_{u2} = (1 - a) Q_u$ ,  $Q_f = K_f S_{f3}$  and  $Q_s = K_s S_s$ , where  $E_p$  signifies the potential evapotranspiration,  $c = 10^{-2}$ ,  $\bar{S}_u = S_u / S_{u,max}$  and  $S_{u,max}$ ,  $a$ ,  $b$ ,  $K_s$  and  $K_f$  are unknown parameters.

2740

2741 hymod parameters with their corresponding symbols, units, and lower and upper bounds.

### 2742 H.2 Hydrologic model

2743 The Hydrologic model (Hmodel) is a parsimonious conceptual watershed model originally developed  
 2744 by *Schoups et al.* (2010). This model transforms rainfall into runoff at the watershed outlet using an  
 2745 interception, unsaturated zone, fast and slow flow reservoir, respectively, which simulate interception,

Table H.1: Summary of hmod parameters and their symbols, units, and lower and upper bounds.

Parameter	Symbol	Units	Min.	Max.
Maximum storage unsaturated zone	$S_{u,\max}$	mm	50	1000
Spatial variability of soil moisture capacity	$b$	—	$10^{-1}$	10
Flow partitioning coefficient	$a$	—	0	1
Recession constant, slow reservoir	$K_s$	1/d	$10^{-4}$	1
Recession constant, fast reservoir	$K_f$	1/d	$10^{-1}$	5

throughfall, evaporation, surface runoff, percolation, fast streamflow and baseflow (see Figure H.2). Table

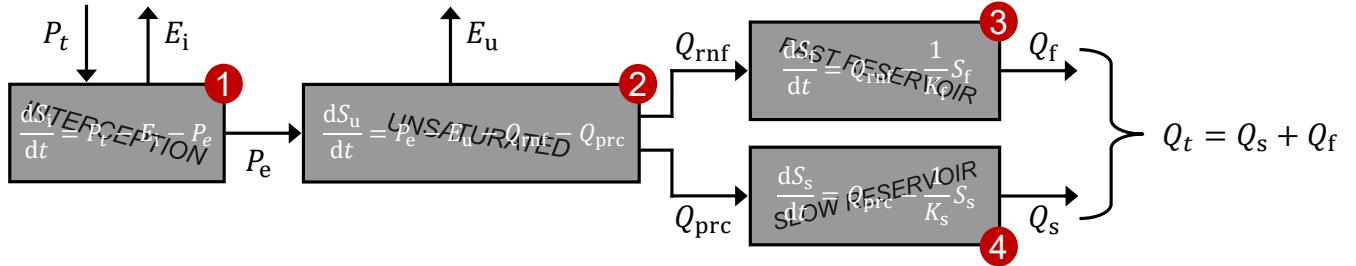


Figure H.2: Schematic illustration of the hmodel after *Schoups et al.* (2010). Grey boxes, labeled in red, correspond to fictitious control volumes of the watershed which govern the rainfall-runoff transformation. The state variables,  $S_i$ ,  $S_u$ ,  $S_f$  and  $S_s$ , correspond to the water storage in each compartment. Arrows portray the fluxes into and out of the compartments, including precipitation,  $P_t$ , interception evaporation,  $E_i$ , excess precipitation,  $P_e$ , evaporation,  $E_u$ , surface runoff,  $Q_{\text{rnf}}$ , percolation,  $Q_{\text{prc}}$ , fast flow,  $Q_f$  and baseflow,  $Q_s$ . The fluxes are computed as follows,  $E_i = E_p f(\bar{S}_i, \alpha_i)$ ,  $P_e = P f(\bar{S}_i, \alpha_p)$ ,  $E_u = (E_p - E_i) f(\bar{S}_u, \alpha_e)$ ,  $Q_{\text{rnf}} = P_e f(\bar{S}_u, \alpha_f)$ ,  $Q_{\text{prc}} = Q_{\max} f(\bar{S}_u, \alpha_f)$ ,  $Q_f = S_f / K_f$  and  $Q_s = S_s / K_s$ , where  $E_p$  signifies the potential evapotranspiration, the functions,  $f(x, y) = (1 - g(-xy)) / (1 - g(-x))$  and  $g(c) = \exp(\min(c, 300))$ , protect against overflow,  $\alpha_i = 50$ ,  $\alpha_p = -50$ ,  $\bar{S}_i = S_i / I_{\max}$ ,  $\bar{S}_u = S_u / S_{\max}$  and  $I_{\max}$ ,  $S_{\max}$ ,  $Q_{\max}$ ,  $\alpha_e$ ,  $\alpha_f$ ,  $K_f$  and  $K_s$  are unknown parameters.

2746

H.2 lists the seven hmodel parameters and their corresponding symbols, units and upper and lower bounds.

Table H.2: Description of Hmodel parameters, including symbols, units, lower and upper bounds.

Parameter	Symbol	Units	Min.	Max.
Maximum interception	$I_{\max}$	mm	0.1	10
Soil water storage capacity	$S_{\max}$	mm	10	1000
Maximum percolation rate	$Q_{\max}$	mm/d	$10^{-1}$	100
Evaporation parameter	$\alpha_e$	—	0	100
Runoff parameter	$\alpha_f$	—	-10	10
Time constant, fast reservoir	$K_f$	d	$10^{-1}$	10
Time constant, slow reservoir	$K_s$	d	1	150

2748

### 2749 H.3 Sacramento Soil Moisture Accounting model

2750 The Sacramento Soil Moisture Accounting (SAC-SMA) model is used by the National Weather Service  
 2751 River Forecast System (NWSRFS) for flood forecasting throughout the United States. The model  
 2752 converts areal average precipitation into streamflow (*Burnash et al.*, 1973). Our implementation follows  
 2753 *Clark et al.* (2008) and is presented in Figure H.3. Table H.3 presents the thirteen SAC-SMA model  
 parameters with their corresponding symbols, units, and lower and upper bounds. This concludes the

Table H.3: Description of SAC-SMA parameters, including symbols, units, lower and upper bounds.

Parameter	Symbol	Units	Min.	Max.
Upper zone tension water maximum storage	$S_{T,\max}^1$	mm	50	500
Upper zone free water maximum storage	$S_{F,\max}^1$	mm	10	500
Lower zone tension water maximum storage	$S_{T,\max}^2$	mm	10	500
Lower zone free water primary maximum storage	$S_{FP,\max}^2$	mm	10	1000
Lower zone free water supplemental maximum storage	$S_{FS,\max}^2$	mm	10	1000
Percolation multiplier for the lower layer	$\alpha$	-	1	250
Percolation exponent for the lower layer	$\psi$	-	1	5
Upper zone free water lateral depletion rate (interflow rate)	$k_i$	mm/d	$10^{-2}$	100
Fraction of percolation to tension storage in the lower layer	$\kappa$	-	0.05	0.95
Base flow depletion rate for primary reservoir	$\nu_P$	$d^{-1}$	$10^{-3}$	0.25
Base flow depletion rate for secondary reservoir	$\nu_S$	$d^{-1}$	$10^{-3}$	0.25
Maximum saturated area (fraction)	$A_{c,\max}$	-	0.05	0.95
Recession constant, routing	$K_f$	$d^{-1}$	$10^{-1}$	5

2754

2755 description of the models.

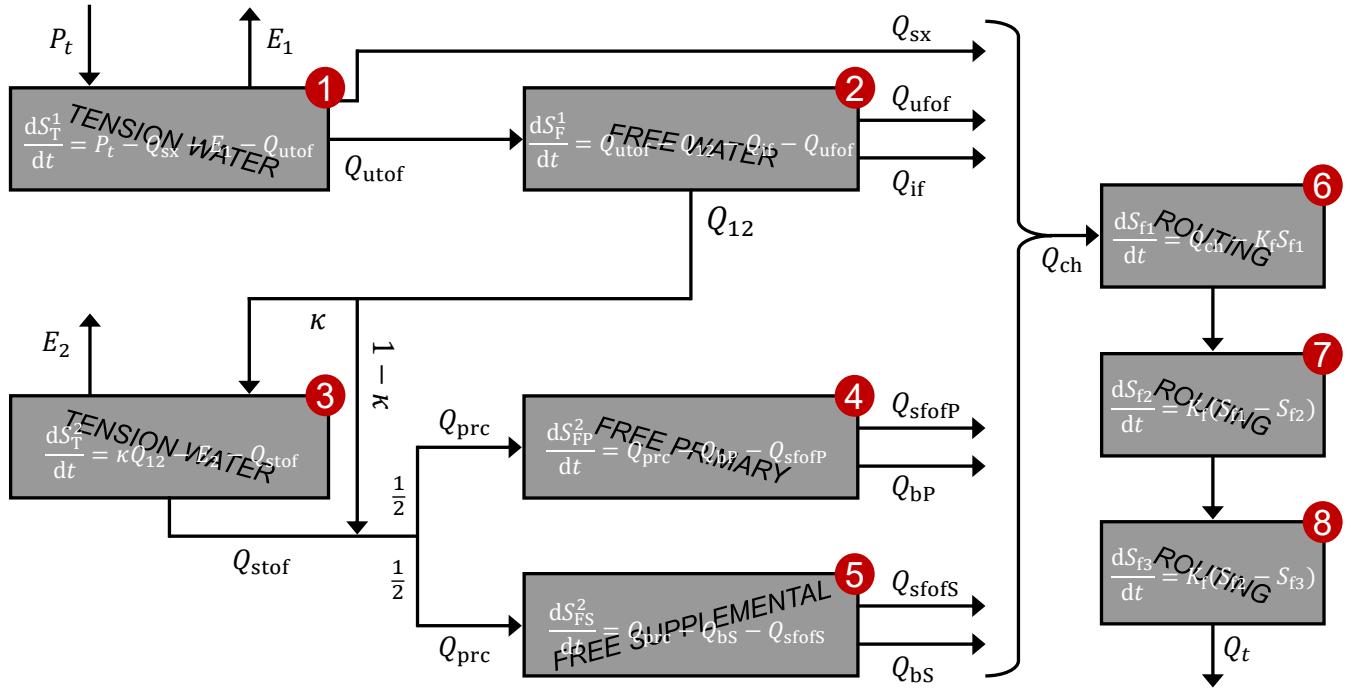


Figure H.3: Schematic illustration of the SAC-SMA model after *Burnash et al. (1973)* and *Clark et al. (2008)*. Grey boxes, labeled in red, correspond to fictitious control volumes of the watershed which govern the rainfall-runoff transformation. The model has eight state variables, including the tension water content in the upper soil layer,  $S_T^1$ , the free water content in the upper soil layer,  $S_F^1$ , the tension water content in the lower soil layer,  $S_T^2$ , the free water content in the primary baseflow reservoir,  $S_{FP}^2$ , the free water content in the secondary baseflow reservoir,  $S_{FS}^2$ , and the water storage,  $S_{f1}$ ,  $S_{f2}$  and  $S_{f3}$ , in the cascade of three linear reservoirs used for routing. Arrows portray the fluxes into and out of the compartments, including precipitation,  $P_t$ , evaporation from upper soil layer,  $E_1$ , overflow of water from tension storage in the upper soil layer,  $Q_{utof}$ , surface runoff,  $Q_{sx}$ , overflow of water from free storage in the upper soil layer,  $Q_{ufof}$ , interflow,  $Q_{if}$ , percolation of water from the upper to the lower layer,  $Q_{12}$ , evaporation from the lower soil layer,  $E_2$ , overflow of water from tension storage in the lower soil layer,  $Q_{stof}$ , water flow into primary and supplemental storage,  $Q_{prc}$ , overflow of water from primary base flow storage in the lower soil layer,  $Q_{sfop}$ , overflow of water from secondary base flow storage in the lower soil layer,  $Q_{sfos}$ , base flow from the primary reservoir,  $Q_{bP}$ , and baseflow from the secondary reservoir,  $Q_{bs}$ . These fluxes are computed as follows,  $E_1 = E_p(S_T^1/S_{T,max}^1)$ ,  $Q_{sx} = A_{c,max}(S_T^1/S_{T,max}^1)P_t$ ,  $Q_{utof} = (P_t - Q_{sx})f(S_T^1, S_{T,max}^1)$ ,  $Q_{12} = Q_0d_{lz}(S_T^1/S_{F,max}^1)$ ,  $Q_{ufof} = Q_{utof}f(S_F^1, S_{F,max}^1)$ ,  $Q_{if} = k_i(S_F^1/S_{F,max}^1)$ ,  $E_2 = (E_p - E_1)(S_T^2/S_{T,max}^2)$ ,  $Q_{stof} = \kappa Q_{12}f(S_T^2, S_{T,max}^2)$ ,  $Q_{prc} = \frac{1}{2}(1 - \kappa)Q_{12} + \frac{1}{2}Q_{stof}$ ,  $Q_{sfop} = Q_{prc}f(S_{FP}^2, S_{FP,max}^2)$ ,  $Q_{sfos} = Q_{prc}f(S_{FS}^2, S_{FS,max}^2)$ ,  $Q_{bP} = \nu_P S_{FP}^2$ ,  $Q_{bs} = \nu_S S_{FS}^2$ , where  $E_p$  signifies the potential evapotranspiration,  $Q_0 = \nu_P S_{FP,max}^2 + \nu_S S_{FS,max}^2$ ,  $d_{lz} = 1 + \alpha[(S_T^2 + S_{FP}^2 + S_{FS}^2)/(S_{T,max}^2 + S_{FP,max}^2 + S_{FS,max}^2)]^\psi$ , the smoothing function  $f(x, y) = 1/(1 + \exp\{-[x - (y - \epsilon\rho y)]/\rho y\})$ , with  $\rho = 10^{-2}$  and  $\epsilon = 5$ , and  $S_{F,max}^1$ ,  $S_{T,max}^1$ ,  $S_{FP,max}^2$ ,  $S_{FS,max}^2$ ,  $S_{T,max}^2$ ,  $\alpha$ ,  $\psi$ ,  $k_i$ ,  $\kappa$ ,  $\nu_P$ ,  $\nu_S$  and  $A_{c,max}$  are unknown parameters. The channel inflow,  $Q_{ch} = Q_{sx} + Q_{ufof} + Q_{if} + Q_{QstofP} + Q_{bP} + Q_{Qstofs} + Q_{bs}$ , is routed through three linear reservoirs with a common recession constant,  $K_f$ , to yield the streamflow,  $Q_t = K_f S_{f3}$ .

## 2756 Appendix I: Continuous Ranked Probability Score of $\mathcal{P}_{\text{III}}(\mu, \sigma^2, \rho)$

2757 The continuous ranked probability score (CRPS) is equal to the integral of the quantile scores

$$2758 S_{\text{CRPS}}(P, \omega) = - \int_{-\infty}^{\infty} S^{\tau} (F_P(z) - \mathbb{1}\{z \geq \omega\})^2 dz, \quad (\text{I.1})$$

2759 where  $F_P(z)$  denotes the cumulative distribution function of  $P$  and the indicator function,  $\mathbb{1}\{a\}$ , returns  
 2760 1 if  $a$  is true and zero otherwise. Suppose that the distribution forecast  $P$  follows a univariate Pearson  
 2761 type III distribution  $\mathcal{P}_{\text{III}}(\mu, \sigma^2, \rho)$  with mean  $\mu$ , variance  $\sigma^2$  and skewness  $\rho$ . If we reparametrize the PIII  
 2762 distribution and define  $\xi = \mu - 2\sigma/\rho$ ,  $a = 4/\rho^2$  and  $b = \frac{1}{2}\sigma|\rho|$  as location, shape and scale parameters,  
 2763 respectively, then the cumulative distribution function of  $P = \mathcal{P}_{\text{III}}(\xi, a, b)$  simplifies to (Tegos *et al.*, 2022)

$$2764 F_P(x, \xi, a, b) = \begin{cases} \frac{1}{\Gamma(a)} \gamma(a, b^{-1}(x - \xi)) & \text{if } \rho > 0 \\ \frac{1}{\Gamma(a)} \Gamma(a, b^{-1}(\xi - x)) & \text{if } \rho < 0 \end{cases} \quad (\text{I.2})$$

2765 where  $x, \xi, a, b \in \mathbb{R}$ ,  $a > 0$ ,  $b > 0$  and

$$2766 \Gamma(a, q) = \int_q^{\infty} t^{a-1} \exp(-t) dt \quad (\text{I.3})$$

2767 and

$$2768 \gamma(a, q) = \int_0^q t^{a-1} \exp(-t) dt \quad (\text{I.4})$$

2769 denote the upper and lower incomplete gamma functions, respectively. If  $\rho > 0$  then  $x \in (\xi, \infty)$  and if  
 2770  $\rho < 0$  then  $x \in (-\infty, \xi)$ . Next, we derive an analytic expression for the CRPS of  $P = \mathcal{P}_{\text{III}}(\xi, a, b)$  and  
 2771 verifying observation  $\omega \in \Omega$ .

### 2772 I.1 Analytic expression for $S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega)_{\rho > 0}$ : positive skewness

2773 Let us first assume that  $\rho > 0$  and, thus,  $x \in [\xi, \infty)$ . In our derivation we work with  $z = x - \xi$  and, thus,  
 2774  $z \in [0, \infty)$ . This change of variables simplifies an analytic solution of the CRPS as we will demonstrate  
 2775 next. Equation (I.1) is now equal to

$$2776 2777 S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega)_{\rho > 0} = - \int_0^{z_{\omega}} \left[ \frac{1}{\Gamma(a)} \gamma(a, b^{-1}z) \right]^2 dz - \int_{z_{\omega}}^{\infty} \left[ \frac{1}{\Gamma(a)} \gamma(a, b^{-1}z) - 1 \right]^2 dz, \quad (\text{I.5})$$

2778 where  $z_{\omega} = \omega - \xi$ . We first derive an expression for the left integral using integration by parts

$$2779 2780 \int u dv = uv - \int v du, \quad (\text{I.6})$$

2781 with  $u = \gamma^2(a, b^{-1}z)$  and  $dv = 1dz$ . Then

$$2782 \quad du = \frac{2}{z} \exp(-b^{-1}z)(b^{-1}z)^a \gamma(a, b^{-1}z) \quad (I.7)$$

2783

2784 and  $v = z$  to yield

$$2785 \quad -\frac{1}{\Gamma^2(a)} \int \gamma^2(a, b^{-1}z) dz = -\frac{1}{\Gamma^2(a)} \left( z \gamma^2(a, b^{-1}z) - \int z \frac{2}{z} \exp(-b^{-1}z)(b^{-1}z)^a \gamma(a, b^{-1}z) \right)$$

2786

$$= -\frac{1}{\Gamma^2(a)} \left( z \gamma^2(a, b^{-1}z) - 2 \int \exp(-b^{-1}z)(b^{-1}z)^a \gamma(a, b^{-1}z) \right). \quad (I.8)$$

2787

2788 The online calculator of Wolfram|Alpha returns a closed-form expression for the indefinite integral

$$2789 \quad \int \exp(-b^{-1}z)(b^{-1}z)^a \gamma(a, b^{-1}z) dz = \int \exp(-b^{-1}z)(b^{-1}z)^a (\Gamma(a) - \Gamma(a, b^{-1}z)) dz$$

2790

$$= b \exp(-b^{-1}z)(b^{-1}z)^a \Gamma(a, b^{-1}z) + \frac{1}{2} ab \Gamma^2(a, b^{-1}z)$$

2791

$$- 4^{-a} b \Gamma(2a, 2b^{-1}z) - b \Gamma(a) \Gamma(a+1, b^{-1}z) + C. \quad (I.9)$$

2792

2793 If we substitute Equation (I.9) into Equation (I.8)

$$2794 \quad -\frac{1}{\Gamma^2(a)} \int \gamma^2(a, z) dz = -\frac{1}{\Gamma^2(a)} \left( z \gamma^2(a, b^{-1}z) - 2(b \exp(-b^{-1}z)(b^{-1}z)^a \Gamma(a, b^{-1}z) + \frac{1}{2} ab \Gamma^2(a, b^{-1}z) \right.$$

2795

$$\left. - 4^{-a} b \Gamma(2a, 2b^{-1}z) - b \Gamma(a) \Gamma(a+1, b^{-1}z)) \right) + C, \quad (I.10)$$

2796

2797 and admit the integral limits

$$2798 \quad -\frac{1}{\Gamma^2(a)} \int_0^{z_\omega} \gamma^2(a, b^{-1}z) dz = -\frac{1}{\Gamma^2(a)} \left| z \gamma^2(a, b^{-1}z) - 2(b \exp(-b^{-1}z)(b^{-1}z)^a \Gamma(a, b^{-1}z) \right.$$

2799

$$\left. + \frac{1}{2} ab \Gamma^2(a, b^{-1}z) - 4^{-a} b \Gamma(2a, 2b^{-1}z) - b \Gamma(a) \Gamma(a+1, b^{-1}z) \right|_0^{z_\omega}, \quad (I.11)$$

2800

2801 then we yield the following expression for the left integral of Equation (I.5)

$$2802 \quad -\frac{1}{\Gamma^2(a)} \int_0^{z_\omega} \gamma^2(a, b^{-1}z) dz = -\frac{1}{\Gamma^2(a)} \left( z_\omega \gamma^2(a, b^{-1}z_\omega) - 2(b \exp(-b^{-1}z_\omega)(b^{-1}z_\omega)^a \Gamma(a, b^{-1}z_\omega) \right.$$

2803

$$\left. + \frac{1}{2} ab \Gamma^2(a, b^{-1}z_\omega) - 4^{-a} b \Gamma(2a, 2b^{-1}z_\omega) - b \Gamma(a) \Gamma(a+1, b^{-1}z_\omega)) \right)$$

2804

$$+ \frac{1}{\Gamma^2(a)} \left( -2 \left( \frac{1}{2} ab \Gamma^2(a) - 4^{-a} b \Gamma(2a) - b \Gamma(a) \Gamma(a+1) \right) \right)$$

2805

$$= -\frac{1}{\Gamma^2(a)} \left( z_\omega \gamma^2(a, b^{-1}z_\omega) - 2(b \exp(-b^{-1}z_\omega)(b^{-1}z_\omega)^a \Gamma(a, b^{-1}z_\omega) \right.$$

2806

$$\left. + \frac{1}{2} ab \Gamma^2(a, b^{-1}z_\omega) - 4^{-a} b \Gamma(2a, 2b^{-1}z_\omega) - b \Gamma(a) \Gamma(a+1, b^{-1}z_\omega)) \right)$$

2807

$$+ \frac{2}{\Gamma^2(a)} \left( 4^{-a} b \Gamma(2a) + b \Gamma(a) \Gamma(a+1) \right) - ab. \quad (I.12)$$

2808

2809 Next, we proceed with the right integral of Equation (I.5)

$$2810 \quad - \int_{z_\omega}^{\infty} \left[ \frac{1}{\Gamma(a)} \gamma(a, b^{-1}z) - 1 \right]^2 dz = - \int_{z_\omega}^{\infty} \left( \frac{1}{\Gamma^2(a)} \gamma^2(a, b^{-1}z) - \frac{2}{\Gamma(a)} \gamma(a, b^{-1}z) + 1 \right) dz$$

2811

$$= -\frac{1}{\Gamma^2(a)} \int_{z_\omega}^{\infty} \gamma^2(a, b^{-1}z) dz + \frac{2}{\Gamma(a)} \int_{z_\omega}^{\infty} \gamma(a, b^{-1}z) dz - \int_{z_\omega}^{\infty} dz. \quad (I.13)$$

2812

2813 The first of three integrals is equal to Equation (I.10), and the other two integrals yield

$$\begin{aligned}
 2814 \quad & \frac{2}{\Gamma(a)} \int_{z_\omega}^{\infty} \gamma(a, b^{-1}z) dz - \int_{z_\omega}^{\infty} dz = \left| \frac{2}{\Gamma(a)} (z\Gamma(a) - z\Gamma(a, b^{-1}z) + b\Gamma(a+1, b^{-1}z)) - z \right|_{z_\omega}^{\infty} \\
 2815 \quad & = \left| z - \frac{2}{\Gamma(a)} (z\Gamma(a, b^{-1}z) - b\Gamma(a+1, b^{-1}z)) \right|_{z_\omega}^{\infty}. \tag{I.14}
 \end{aligned}$$

2817 Thus, Equation (I.13) becomes

$$\begin{aligned}
 2818 \quad & - \int_{z_\omega}^{\infty} \left[ \frac{1}{\Gamma(a)} \gamma(a, b^{-1}z) - 1 \right]^2 dz = - \frac{1}{\Gamma^2(a)} \left| z\gamma^2(a, b^{-1}z) - 2(b \exp(-b^{-1}z)(b^{-1}z)^a \Gamma(a, b^{-1}z) \right. \\
 2819 \quad & \quad \left. + \frac{1}{2} ab\Gamma^2(a, b^{-1}z) - 4^{-a} b\Gamma(2a, 2b^{-1}z) - b\Gamma(a)\Gamma(a+1, b^{-1}z)) \right|_{z_\omega}^{\infty} \\
 2820 \quad & \quad + \left| z - \frac{2}{\Gamma(a)} (z\Gamma(a, b^{-1}z) - b\Gamma(a+1, b^{-1}z)) \right|_{z_\omega}^{\infty} \\
 2821 \quad & = - \frac{1}{\Gamma^2(a)} \left| z\gamma^2(a, b^{-1}z) - 2(b \exp(-b^{-1}z)(b^{-1}z)^a \Gamma(a, b^{-1}z) \right. \\
 2822 \quad & \quad \left. + \frac{1}{2} ab\Gamma^2(a, b^{-1}z) - 4^{-a} b\Gamma(2a, 2b^{-1}z) - b\Gamma(a)\Gamma(a+1, b^{-1}z)) - z\Gamma^2(a) \right. \\
 2823 \quad & \quad \left. + 2\Gamma(a)(z\Gamma(a, b^{-1}z) - b\Gamma(a+1, b^{-1}z)) \right|_{z_\omega}^{\infty} \\
 2824 \quad & = \lim_{z \rightarrow \infty} - \frac{1}{\Gamma^2(a)} \left( z \cancel{\gamma^2(a, b^{-1}z)}^{\Gamma^2(a)} - 2(b \exp(-b^{-1}z)(b^{-1}z)^a \Gamma(a, b^{-1}z) \right. \\
 2825 \quad & \quad \left. + \frac{1}{2} ab\Gamma^2(a, b^{-1}z) \cancel{- 4^{-a} b\Gamma(2a, 2b^{-1}z)}^0 - b\Gamma(a)\Gamma(a+1, b^{-1}z) \cancel{- z\Gamma^2(a)}^0 \right. \\
 2826 \quad & \quad \left. + 2\Gamma(a)(z\Gamma(a, b^{-1}z) \cancel{- b\Gamma(a+1, b^{-1}z)}^0) \right) + \frac{1}{\Gamma^2(a)} \left( z_\omega \gamma^2(a, b^{-1}z_\omega) \right. \\
 2827 \quad & \quad \left. - 2(b \exp(-b^{-1}z_\omega)(b^{-1}z_\omega)^a \Gamma(a, b^{-1}z_\omega) + \frac{1}{2} ab\Gamma^2(a, b^{-1}z_\omega) \right. \\
 2828 \quad & \quad \left. - 4^{-a} b\Gamma(2a, 2b^{-1}z_\omega) - b\Gamma(a)\Gamma(a+1, b^{-1}z_\omega)) - z_\omega \Gamma^2(a) \right. \\
 2829 \quad & \quad \left. + 2\Gamma(a)(z_\omega \Gamma(a, b^{-1}z_\omega) - b\Gamma(a+1, b^{-1}z_\omega)) \right) \\
 2830 \quad & = \frac{1}{\Gamma^2(a)} \left( z_\omega \gamma^2(a, b^{-1}z_\omega) - 2(b \exp(-b^{-1}z_\omega)(b^{-1}z_\omega)^a \Gamma(a, b^{-1}z_\omega) \right. \\
 2831 \quad & \quad \left. + \frac{1}{2} ab\Gamma^2(a, b^{-1}z_\omega) - 4^{-a} b\Gamma(2a, 2b^{-1}z_\omega) - b\Gamma(a)\Gamma(a+1, b^{-1}z_\omega)) \right. \\
 2832 \quad & \quad \left. - z_\omega \Gamma^2(a) + 2\Gamma(a)(z_\omega \Gamma(a, b^{-1}z_\omega) - b\Gamma(a+1, b^{-1}z_\omega)) \right). \tag{I.15}
 \end{aligned}$$

2834 The sum of Equations (I.12) and (I.15) equals  $S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega)_{\rho>0}$ , whence we can write

$$\begin{aligned}
 2835 \quad & S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega)_{\rho>0} = - \frac{1}{\Gamma^2(a)} \left( z_\omega \gamma^2(a, b^{-1}z_\omega) - 2(b \exp(-b^{-1}z_\omega)(b^{-1}z_\omega)^a \Gamma(a, b^{-1}z_\omega) \right. \\
 2836 \quad & \quad \left. + \frac{1}{2} ab\Gamma^2(a, b^{-1}z_\omega) - 4^{-a} b\Gamma(2a, 2b^{-1}z_\omega) - b\Gamma(a)\Gamma(a+1, b^{-1}z_\omega)) \right) \\
 2837 \quad & \quad + \frac{2}{\Gamma^2(a)} \left( 4^{-a} b\Gamma(2a) + b\Gamma(a)\Gamma(a+1) \right) - ab
 \end{aligned}$$

$$\begin{aligned}
& + \frac{1}{\Gamma^2(a)} \left( z_\omega \gamma^2(a, b^{-1}z_\omega) - 2(b \exp(-b^{-1}z_\omega)(b^{-1}z_\omega)^a \Gamma(a, b^{-1}z_\omega) \right. \\
& + \frac{1}{2} ab \Gamma^2(a, b^{-1}z_\omega) - 4^{-a} b \Gamma(2a, 2b^{-1}z_\omega) - b \Gamma(a) \Gamma(a+1, b^{-1}z_\omega) \left. \right) \\
& - z_\omega \Gamma^2(a) + 2\Gamma(a) (z_\omega \Gamma(a, b^{-1}z_\omega) - b \Gamma(a+1, b^{-1}z_\omega)) \Big) \\
& = \frac{2}{\Gamma^2(a)} \left( 4^{-a} b \Gamma(2a) + b \Gamma(a) \Gamma(a+1) \right) - ab - \frac{1}{\Gamma^2(a)} \left( z_\omega \Gamma^2(a) \right. \\
& \left. - 2\Gamma(a) (z_\omega \Gamma(a, b^{-1}z_\omega) - b \Gamma(a+1, b^{-1}z_\omega)) \right) \\
& = 2 \frac{4^{-a} b \Gamma(2a)}{\Gamma^2(a)} + 2b \frac{\Gamma(a+1)}{\Gamma(a)} - ab - z_\omega + 2 \frac{z_\omega \Gamma(a, b^{-1}z_\omega)}{\Gamma(a)} \\
& \quad - 2b \frac{\Gamma(a+1, b^{-1}z_\omega)}{\Gamma(a)} \\
& = 2 \frac{4^{-a} b}{B(a, a)} + ab - z_\omega + 2z_\omega (1 - F_G(z_\omega, a, b)) \\
& \quad - 2ab (1 - F_G(z_\omega, a+1, b)), \tag{I.16}
\end{aligned}$$

where  $B(u, v) = \Gamma(u)\Gamma(v)/\Gamma(u+v)$  is the beta function of the first kind and  $F_G(z, a, b)$  is the CDF of the gamma distribution

$$F_G(z, a, b) = \frac{1}{\Gamma(a)} \gamma(a, b^{-1}z), \tag{I.17}$$

with  $z > 0$ , unitless shape parameter  $a > 0$  and scale parameter  $b > 0$ . We can rearrange Equation (I.16) to yield a final expression for the CRPS of a PIII distribution forecast  $P = \mathcal{P}_{\text{III}}(\xi, a, b)$  with positive skewness,  $\rho > 0$  and  $z_\omega = \omega - \xi$

$$S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega)_{\rho > 0} = 2 \frac{4^{-a} b}{B(a, a)} + ab (2F_G(z_\omega, a+1, b) - 1) + z_\omega (1 - 2F_G(z_\omega, a, b)). \tag{I.18}$$

## I.2 Analytic expression for $S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega)_{\rho < 0}$ : negative skewness

Let us now assume that  $\rho < 0$  and, thus,  $x \in (-\infty, \xi]$ . We define  $z = \xi - x$  and, thus,  $z \in [0, \infty)$ . Equation (I.1) becomes

$$S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega)_{\rho < 0} = - \int_0^{z_\omega} \left[ \frac{1}{\Gamma(a)} \Gamma(a, b^{-1}z) \right]^2 dz - \int_{z_\omega}^\infty \left[ \frac{1}{\Gamma(a)} \Gamma(a, b^{-1}z) - 1 \right]^2 dz. \tag{I.19}$$

We take advantage of the following identity

$$\gamma(a, b^{-1}z) + \Gamma(a, b^{-1}z) = \Gamma(a), \tag{I.20}$$

2862 to write Equation (I.19) in another form

$$\begin{aligned}
 2863 \quad S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega)_{\rho < 0} &= - \int_0^{z_\omega} \left[ \frac{1}{\Gamma(a)} (\Gamma(a) - \gamma(a, b^{-1}z)) \right]^2 dz \\
 2864 \quad &\quad - \int_{z_\omega}^{\infty} \left[ \frac{1}{\Gamma(a)} (\Gamma(a) - \gamma(a, b^{-1}z)) - 1 \right]^2 dz \\
 2865 \quad &= - \int_0^{z_\omega} \left[ \frac{1}{\Gamma(a)} \gamma(a, b^{-1}z) - 1 \right]^2 dz - \int_{z_\omega}^{\infty} \left[ \frac{1}{\Gamma(a)} \gamma(a, b^{-1}z) \right]^2 dz, \quad (I.21)
 2866
 \end{aligned}$$

2867 which is identical to Equation (I.5) but with left and right integrals swapped. This confirms that

$$2868 \quad S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega)_{\rho < 0} = S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega)_{\rho > 0}, \quad (I.22)$$

2870 and, thus, we yield the following expression for the CRPS of a PIII distribution forecast  $P = \mathcal{P}_{\text{III}}(\xi, a, b)$

2871 with  $\rho \neq 0$  and  $z_\omega = \xi - \omega$

$$\begin{aligned}
 2872 \quad S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega)_{\rho < 0} &= 2 \frac{4^{-a}b}{B(a, a)} + ab(2F_G(z_\omega, a+1, b) - 1) + z_\omega(1 - 2F_G(z_\omega, a, b)), \quad (I.23) \\
 2873
 \end{aligned}$$

### 2874 I.3 Analytic expression for $S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega)$ : positive/negative skewness

2875 We can combine the mathematical expressions of  $S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega)_{\rho > 0}$  and  $S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega)_{\rho < 0}$   
 2876 into one general Equation for the CRPS of  $P = \mathcal{P}_{\text{III}}(\xi, a, b)$  and verifying observation  $\omega \in \Omega$

$$\begin{aligned}
 2877 \quad S_{\text{CRPS}}(\mathcal{P}_{\text{III}}(\xi, a, b), \omega) &= 2 \frac{4^{-a}b}{B(a, a)} - ab + |\omega - \xi| + 2abF_G(|\omega - \xi|, a+1, b) \\
 2878 \quad &\quad - 2|\omega - \xi|F_G(|\omega - \xi|, a, b), \quad (I.24)
 \end{aligned}$$

2880 where  $z_\omega = |\omega - \xi|$ . The first term in the above expression is equal to the first term,  $\frac{1}{2}\mathbb{E}_P[|y - y^*|]$ , of  
 2881 Equation (63) and the sum of all the remaining terms of Equation (I.24) equals  $-\mathbb{E}_P[|y - \omega|]$ .

2882 We would be remiss not to address two well-known limiting cases of the PIII distribution. For  $\xi = 0$ ,  
 2883 thus,  $\mu = 2\sigma/\rho$ , the PIII distribution  $\mathcal{P}_{\text{III}}(\xi, a, b)$  reduces to the gamma distribution  $\mathcal{G}(a, b)$  and Equation  
 2884 (I.24) simplifies to

$$\begin{aligned}
 2885 \quad S_{\text{CRPS}}(\mathcal{G}(a, b), \omega) &= 2 \frac{4^{-a}b}{B(a, a)} - ab + |\omega| + 2abF_G(|\omega|, a+1, b) - 2|\omega|F_G(|\omega|, a, b). \quad (I.25) \\
 2886
 \end{aligned}$$

2887 This expression for the CRPS of  $P = \mathcal{G}(a, b)$  matches the numerical estimates of the CRPS shown in  
 2888 Figure 10 using the solid yellow line. If  $\rho = 0$ , then the PIII distribution  $\mathcal{P}_{\text{III}}(\mu, \sigma^2, 0)$  simplifies to a  
 2889 normal distribution  $\mathcal{N}(\mu, \sigma^2)$  and the CRPS can be computed using Equation (G.12). This concludes  
 2890 the derivation.

## 2891 Appendix J: Continuous Ranked Probability Score of $\mathcal{GEV}(\mu, \sigma^2, \xi)$

2892 We revisit the quantile form of the continuous ranked probability score (CRPS)

$$2893 S_{\text{CRPS}}(P, \omega) = \omega(1 - 2F_P(\omega)) + 2 \int_0^1 \tau F_P^{-1}(\tau) d\tau - 2 \int_{F_P(\omega)}^1 F_P^{-1}(\tau) d\tau. \quad (\text{J.1})$$

2895 Suppose that the distribution forecast  $P$  follows a generalized extreme value distribution  $\mathcal{GEV}(\mu, \sigma^2, \xi)$

2896 with mean  $\mu \in \mathbb{R}$ , variance  $\sigma^2 > 0$  and shape parameter  $\xi \in \mathbb{R}$ . The CDF of the GEV distribution equals

$$2897 F_{\mathcal{GEV}}(x, \mu, \sigma^2, \xi) = \begin{cases} \exp\left[-\exp\left(-\frac{\xi}{\sigma}(x - \mu)\right)\right] & \text{if } \xi = 0 \\ \exp\left[-\left(1 + \frac{\xi}{\sigma}(x - \mu)\right)^{-1/\xi}\right] & \text{if } \xi < 0 \text{ and } x < -\frac{1}{\xi} \\ 1 & \text{if } \xi < 0 \text{ and } x \geq -\frac{1}{\xi} \\ 0 & \text{if } \xi > 0 \text{ and } x \leq -\frac{1}{\xi} \\ \exp\left[-\left(1 + \frac{\xi}{\sigma}(x - \mu)\right)^{-1/\xi}\right] & \text{if } \xi > 0 \text{ and } x > -\frac{1}{\xi}, \end{cases} \quad (\text{J.2})$$

2898 and its quantile function has the following explicit expression

$$2899 F_{\mathcal{GEV}}^{-1}(x, \mu, \sigma^2, \xi) = \begin{cases} \mu - \sigma \log_e(-\log_e(\tau)) & \text{if } \xi = 0 \text{ and } \tau \in (0, 1) \\ \mu + \frac{\sigma}{\xi} \left( (-\log_e(\tau))^{-\xi} - 1 \right) & \text{if } \xi \neq 0, \end{cases} \quad (\text{J.3})$$

2900 where  $\tau \in [0, 1]$  if  $\xi > 0$  and  $\tau \in (0, 1]$  if  $\xi < 0$ . Next, we derive an analytic expression for the CRPS  
2901 of  $P = \mathcal{GEV}(\mu, \sigma^2, \xi)$  and verifying observation  $\omega \in \Omega$ . For  $\xi \geq 1$  the CRPS of  $P = \mathcal{GEV}(\mu, \sigma^2, \xi)$  is  
2902 undefined. Thus, we restrict our attention to  $\xi < 1$ .

### 2903 J.1 Analytic expression for $S_{\text{CRPS}}(\mathcal{GEV}(\mu, \sigma^2, \xi), \omega)_{\xi < 1, \xi \neq 0}$

2904 We first consider the case of a non-zero shape parameter,  $\xi \neq 0$ , of the GEV distribution. The indefinite  
2905 form of the first integral of the CRPS in Equation (J.1) can be expressed analytically using integration  
2906 by parts

$$\begin{aligned} 2907 \int \tau F_P^{-1}(\tau) d\tau &= \frac{\tau}{\xi} \left( \mu \xi \tau + \sigma \Gamma(1 - \xi, -\log_e(\tau)) - \sigma \tau \right) + C \\ 2908 &\quad - \int \frac{1}{\xi} \left( \mu \xi \tau + \sigma \Gamma(1 - \xi, -\log_e(\tau)) - \sigma \tau \right) d\tau \\ 2909 &= \frac{\tau}{\xi} \left( \mu \xi \tau + \sigma \Gamma(1 - \xi, -\log_e(\tau)) - \sigma \tau \right) \\ 2910 &\quad - \frac{1}{2\xi} \left[ \tau \left( \mu \tau \xi + 2\sigma \Gamma(1 - \xi, -\log_e(\tau)) - \sigma \tau \right) - 2^\xi \sigma \Gamma(1 - \xi, -2 \log_e(\tau)) \right] + C, \quad (\text{J.4}) \end{aligned}$$

2912 where

$$2913 \quad \Gamma(a, q) = \int_q^\infty t^{a-1} \exp(-t) dt \quad (\text{J.5})$$

2914 denotes the upper incomplete gamma function. If we now admit the lower and upper limits of the  
2915 quantiles, we yield

$$\begin{aligned} 2916 \quad \int_0^1 \tau F_P^{-1}(\tau) d\tau &= \left[ \frac{\tau}{\xi} \left( \mu \xi \tau + \sigma \Gamma(1 - \xi, -\log_e(\tau)) - \sigma \tau \right) \right. \\ 2917 &\quad \left. - \frac{1}{2\xi} \left[ \tau \left( \mu \tau \xi + 2\sigma \Gamma(1 - \xi, -\log_e(\tau)) - \sigma \tau \right) - 2^\xi \sigma \Gamma(1 - \xi, -2 \log_e(\tau)) \right] + C \right]_0^1 \\ 2918 &= \left( \mu + \frac{\sigma}{\xi} \Gamma(1 - \xi) - \frac{\sigma}{\xi} - \frac{\mu}{2} - \frac{\sigma}{\xi} \Gamma(1 - \xi) + \frac{\sigma}{2\xi} + \frac{2^\xi \sigma \Gamma(1 - \xi)}{2\xi} \right) - (0) \\ 2919 &= \frac{\mu}{2} + \frac{\sigma}{2\xi} (2^\xi \Gamma(1 - \xi) - 1). \end{aligned} \quad (\text{J.6})$$

2921 The second or right integral of the CRPS in Equation (J.1) results in

$$\begin{aligned} 2922 \quad \int_{F_P(\omega)}^1 F_P^{-1}(\tau) d\tau &= \left[ \frac{1}{\xi} \left( \mu \xi \tau + \sigma \Gamma(1 - \xi, -\log_e(\tau)) - \sigma \tau \right) + C \right]_{F_P(\omega)}^1 \\ 2923 &= \left( \mu + \frac{\sigma}{\xi} \Gamma(1 - \xi) - \frac{\sigma}{\xi} \right) - \left( \mu F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, \xi) \right. \\ 2924 &\quad \left. + \frac{\sigma}{\xi} \Gamma(1 - \xi, -\log_e(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, \xi))) \right) - \frac{\sigma}{\xi} F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, \xi) \\ 2925 &= \mu - \mu F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, \xi) + \frac{\sigma}{\xi} \left[ \Gamma(1 - \xi) + F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, \xi) \right. \\ 2926 &\quad \left. - \Gamma(1 - \xi, -\log_e(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, \xi))) \right] - 1. \end{aligned} \quad (\text{J.7})$$

2928 Next, we can insert the analytic expressions of the two integrals into Equation (J.1)

$$\begin{aligned} 2929 \quad S_{\text{CRPS}}(\mathcal{GEV}(\omega, \mu, \sigma^2, \xi), \omega)_{\xi < 1, \xi \neq 0} &= \omega (1 - 2F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, \xi)) + \mu + \frac{\sigma}{\xi} (2^\xi \Gamma(1 - \xi) - 1) \\ 2930 &\quad - 2\mu + 2\mu F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, \xi) - \frac{2\sigma}{\xi} \left[ \Gamma(1 - \xi) + F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, \xi) \right. \\ 2931 &\quad \left. - \Gamma(1 - \xi, -\log_e(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, \xi))) \right] - 1. \end{aligned}$$

2933 The above expression may be rearranged and rewritten to yield

$$\begin{aligned} 2934 \quad S_{\text{CRPS}}(\mathcal{GEV}(\omega, \mu, \sigma^2, \xi), \omega)_{\xi < 1, \xi \neq 0} &= (\mu - \omega) (2F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, \xi) - 1) + \frac{\sigma}{\xi} \left( 1 - (2 - 2^\xi) \Gamma(1 - \xi) \right) \\ 2935 &\quad + \frac{2\sigma}{\xi} \left[ \Gamma(1 - \xi, -\log_e(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, \xi))) \right. \\ 2936 &\quad \left. - F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, \xi) \right]. \end{aligned} \quad (\text{J.8})$$

## 2938 J.2 Analytic expression for $S_{\text{CRPS}}(\mathcal{GEV}(\mu, \sigma^2, 0), \omega)$

2939 Next, we consider  $\xi = 0$ . We use integration by parts

$$2940 \quad \int u dv = uv - \int v du, \quad (\text{J.9})$$

2942 and yield the following expression for the first of two integrals of Equation (J.1)

$$\begin{aligned} 2943 \quad \int \tau F_P^{-1}(\tau) d\tau &= \frac{1}{2}\mu\tau^2 - \frac{1}{2}\sigma\tau^2 \log_e(-\log_e(\tau)) + C - \int -\frac{\sigma\tau^2}{2\tau \log_e(\tau)} d\tau \\ 2944 \quad &= \frac{1}{2}\mu\tau^2 - \frac{1}{2}\sigma\tau^2 \log_e(-\log_e(\tau)) + C + \frac{1}{2}\sigma \int \frac{\tau}{\log_e(\tau)} d\tau \\ 2945 \quad &= \frac{1}{2}\mu\tau^2 - \frac{1}{2}\sigma\tau^2 \log_e(-\log_e(\tau)) + \frac{1}{2}\sigma \text{Ei}(2\log_e(\tau)) + C, \end{aligned} \quad (\text{J.10})$$

2947 where  $\text{Ei}(x)$  is the exponential integral function. If we admit the integral limits, we yield

$$\begin{aligned} 2948 \quad \int_0^1 \tau F_P^{-1}(\tau) d\tau &= \left[ \frac{1}{2}\mu\tau^2 - \frac{1}{2}\sigma\tau^2 \log_e(-\log_e(\tau)) + \frac{1}{2}\sigma \text{Ei}(2\log_e(\tau)) + C \right]_0^1 \\ 2949 \quad &= \frac{1}{2}\mu - \frac{1}{2}\sigma \lim_{\tau \rightarrow 1^-} \left( \tau^2 \log_e(-\log_e(\tau)) - \text{Ei}(2\log_e(\tau)) \right) \\ 2950 \quad &= \frac{1}{2}\mu + \frac{1}{2}\sigma(\gamma_c + \log_e(2)), \end{aligned} \quad (\text{J.11})$$

2952 where  $\gamma_c = 0.57721566\dots$  is the Euler-Mascheroni constant. The right integral of Equation (J.1) becomes

$$\begin{aligned} 2953 \quad \int_{F_P(\omega)}^1 F_P^{-1}(\tau) d\tau &= \left[ \mu\tau + \sigma \text{Li}(\tau) - \sigma\tau \log_e(-\log_e(\tau)) \right]_{F_P(\omega)}^1 \\ 2954 \quad &= \left( \mu + \lim_{\tau \rightarrow 1^-} \left( \sigma \text{Li}(\tau) - \sigma\tau \log_e(-\log_e(\tau)) \right) \right) \\ 2955 \quad &\quad - \left( \mu F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0) + \sigma \text{Li}(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0)) \right. \\ 2956 \quad &\quad \left. - \sigma F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0) \log_e(-\log_e(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0))) \right) \\ 2957 \quad &= \mu + \gamma_c \sigma - \mu F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0) - \sigma \text{Li}(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0)) \\ 2958 \quad &\quad + \sigma F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0) \log_e(-\log_e(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0))). \end{aligned} \quad (\text{J.12})$$

2960 where  $\text{Li}(x)$  signifies the logarithmic integral function. If we substitute Equations (J.11) and (J.12) into  
2961 Equation (J.1), we yield the following expression for the CRPS of  $P = \mathcal{GEV}(\mu, \sigma^2, 0)$

$$\begin{aligned} 2962 \quad S_{\text{CRPS}}(\mathcal{GEV}(\omega, \mu, \sigma^2, 0), \omega) &= \omega(1 - 2F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0)) + \mu + \sigma(\gamma_c + \log_e(2)) - 2\mu \\ 2963 \quad &\quad - 2\gamma_c \sigma + 2\mu F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0) + 2\sigma \text{Li}(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0)) \\ 2964 \quad &\quad - 2\sigma F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0) \log_e(-\log_e(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0))) \\ 2965 \quad &= \omega(1 - 2F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0)) - \mu - \gamma_c \sigma + \sigma \log_e(2) \end{aligned}$$

$$\begin{aligned}
& + 2\mu F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0) + 2\sigma \text{Li}(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0)) \\
& - 2\sigma F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0) \log_e(-\log_e(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0))) \\
& = \omega - \mu - \gamma_c \sigma + \sigma \log_e(2) + 2\sigma \text{Li}(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0)) \\
& - 2[\omega - \mu + \sigma \log_e(-\log_e(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0)))] F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0). \quad (\text{J.13})
\end{aligned}$$

2966 Per the quantile function in Equation (J.3), we find that  
2967  
2968

$$- \mu + \sigma \log_e(-\log_e(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0))) = -\omega, \quad (\text{J.14})$$

2971 and, thus, the last term amounts to zero. This results in the following expression for the CRPS of a  
2972 GEV distribution forecast  $P = \mathcal{GEV}(\mu, \sigma^2, 0)$  with shape parameter  $\xi = 0$   
2973

$$S_{\text{CRPS}}(\mathcal{GEV}(\omega, \mu, \sigma^2, 0), \omega) = \omega - \mu + \sigma \log_e(2) + 2\sigma \text{Li}(F_{\mathcal{GEV}}(\omega, \mu, \sigma^2, 0)) - \gamma_c \sigma. \quad (\text{J.15})$$

2974 This concludes the derivation.  
2975  
2976