



**HAL**  
open science

# A generalised approach for identifying influential data in hydrological modelling

D. Wright, M. Thyer, S. Westra, Benjamin Renard, D. Mcinerney

## ► To cite this version:

D. Wright, M. Thyer, S. Westra, Benjamin Renard, D. Mcinerney. A generalised approach for identifying influential data in hydrological modelling. *Environmental Modelling and Software*, 2019, 111, pp.231-247. 10.1016/j.envsoft.2018.03.004 . hal-02608443

**HAL Id: hal-02608443**

**<https://hal.inrae.fr/hal-02608443v1>**

Submitted on 12 Sep 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# ACCEPTED VERSION

David P. Wright, Mark Thyer, Seth Westra, Benjamin Renard, David McInerney  
**A generalised approach for identifying influential data in hydrological modelling**  
Environmental Modelling and Software, 2019; 111:231-247

© 2018 Elsevier Ltd. All rights reserved.

This manuscript version is made available under the CC-BY-NC-ND 4.0 license  
<http://creativecommons.org/licenses/by-nc-nd/4.0/>

Final publication at <http://dx.doi.org/10.1016/j.envsoft.2018.03.004>

## PERMISSIONS

<https://www.elsevier.com/about/our-business/policies/sharing>

Accepted Manuscript

Authors can share their [accepted manuscript](#):

### Immediately

- via their non-commercial personal homepage or blog
- by updating a [preprint](#) in arXiv or RePEc with the [accepted manuscript](#)
- via their research institute or institutional repository for internal institutional uses or as part of an invitation-only research collaboration work-group
- directly by providing copies to their students or to research collaborators for their personal use
- for private scholarly sharing as part of an invitation-only work group on [commercial sites with which Elsevier has an agreement](#)

### After the embargo period

- via non-commercial hosting platforms such as their institutional repository
- via commercial sites with which Elsevier has an agreement

In all cases [accepted manuscripts](#) should:

- link to the formal publication via its DOI
- bear a CC-BY-NC-ND license – this is easy to do
- if aggregated with other manuscripts, for example in a repository or other site, be shared in alignment with our [hosting policy](#)
- not be added to or enhanced in any way to appear more like, or to substitute for, the published journal article

**22 March 2021**

<http://hdl.handle.net/2440/124123>

1 **A generalised approach for identifying influential data in**  
2 **hydrological modelling**

3  
4  
5 **Authors: David P. Wright<sup>1</sup>, Mark Thyer<sup>1</sup>, Seth Westra<sup>1</sup>, Benjamin Renard<sup>2</sup>, David McInerney<sup>1</sup>**

6 1. School of Civil, Environmental and Mining Engineering, University of Adelaide, Adelaide, 5005, Australia

7 2. Irstea, UR Riverly, Lyon-Villeurbanne center, 69625 Villeurbanne, France

8 **Corresponding Author: David P. Wright, [david.p.wright@adelaide.edu.au](mailto:david.p.wright@adelaide.edu.au)**

9 **Submission to Environmental Modelling & Software**  

---

10  
11 **Highlights:**

- 12 1. Influential data points have a disproportionate impact on model predictions  
13 2. A new generalised Cook's distance accurately identifies influential data points  
14 3. More efficient (<1% computational cost) than standard case-deletion approaches  
15 4. Applies to nonlinear regression and hydrological models with heteroscedastic errors  
16 5. Can be used in a Bayesian framework with priors or data uncertainty

19 **Abstract:**

20 Influence diagnostics are used to identify data points that have a disproportionate impact on model  
21 parameters, performance and/or predictions, providing valuable information for use in model  
22 calibration. Regression-theory influence diagnostics identify influential data by combining the  
23 leverage and the standardised residuals, and are computationally more efficient than case-deletion  
24 approaches. This study evaluates the performance of a range of regression-theory influence  
25 diagnostics on ten case studies with a variety of model structures and inference scenarios including:  
26 nonlinear model response, heteroscedastic residual errors, data uncertainty and Bayesian priors. A  
27 new technique is developed, generalised Cook's distance, that is able to accurately identify the same  
28 influential data as standard case deletion approaches (Spearman rank correlation: 0.93-1.00) at a  
29 fraction of the computational cost (<1%). This is because generalised Cook's distance uses a  
30 generalised leverage formulation which outperforms linear and nonlinear leverage formulations  
31 because it has less restrictive assumptions. Generalised Cook's distance has the potential to enable  
32 influential data to be efficiently identified on a wide variety of hydrological and environmental  
33 modelling problems.

34 **Keywords:** *hydrologic model calibration, influence diagnostics, Cook's distance, generalised leverage*

35

## 36 1. Introduction

37 Hydrological model calibration is a critical component of model development as parameters generally  
38 cannot be determined directly from measurements but are instead inferred indirectly by calibrating  
39 the hydrological model to observed hydrological responses (e.g. daily streamflow) [Beven, 2011].  
40 Studies increasingly have called for the use of “influence diagnostics” [e.g., Foglia *et al.*, 2009; Foglia  
41 *et al.*, 2007; Hill *et al.*, 2015; Wright *et al.*, 2015] to understand the extent to which model calibration  
42 outcomes are determined by a small number of data points that may be erroneous or  
43 unrepresentative of overall catchment behaviour. For example, Wright *et al.* [2015] showed that  
44 removing a single value of daily streamflow from a two-year calibration period could change the  
45 predicted streamflow by more than 25% in a semi-arid catchment. There are a range of influence  
46 diagnostics in the literature that have been used to identify which points are influential; the goal of  
47 this paper is to evaluate a generalised approach to identifying influential points that is both accurate  
48 and computationally efficient.

49 Influence diagnostics can be categorised into two different classes: “case-deletion” influence  
50 diagnostics and “regression-theory” influence diagnostics (see Figure 1). Case-deletion influence  
51 diagnostics measure the influence by censoring (“deleting”) a data point (“case”) from the set of  
52 calibration points, then re-calibrating the model. Once case-deletion has been performed, several  
53 approaches can be used to measure influence. The first approach is to evaluate Cook’s distance [Cook,  
54 1977], which is a commonly used measure of influence [Cook, 1977] and has been used in a large  
55 variety of regression problems [Fox and Weisberg, 2011]. The second approach is to quantify the  
56 difference between original and re-calibrated model parameters, model performance (such as  
57 objective function displacement) and/or model predictions of interest [Wright *et al.*, 2015]. Two  
58 further approaches to measure influence are DFFITS and DFBETA [see Cook and Weisberg, 1982].  
59 These are not considered further in this study because DFFITS is conceptually identical to Cook’s  
60 distance (see Cook and Weisberg [1982]), and DFBETA describes the impact of influential data on

61 individual model parameter estimates only [Fox and Weisberg, 2011], whereas Cook's distance has  
62 the flexibility to look at the impact of influential points on parameters (including their interactions)  
63 and predictions.

64 The case-deletion influence diagnostics are classified as "exact" because they make no assumptions  
65 regarding the type of regression model (linear/nonlinear) or the complexity of the residual error model  
66 (Gaussian, heteroscedastic, autocorrelated etc. - see *McInerney et al.* [2017]). This makes them  
67 particularly attractive for hydrological applications, where the hydrological models are generally  
68 nonlinear and assumptions related to the behaviour of the residuals, such as Gaussianity and  
69 homoscedasticity, are typically not supported by the data. The drawback with case-deletion based  
70 influence diagnostics is the high computational demand associated with re-estimating the parameters  
71 for every data point in the observed data (e.g. for a decade of daily data case-deletion requires ~3650  
72 model re-calibrations). This renders influence analysis using case-deletion potentially infeasible for  
73 anything but the simplest hydrological models. A secondary issue with the case-deletion class is that  
74 anomalous results may arise when calibrating to complex response surfaces with multiple local optima  
75 [*Duan et al.*, 1992; *Kavetski et al.*, 2006], as each re-calibration may lead to parameter sets in different  
76 local optima. This may cause the case-deletion calibrated parameter sets to be different from each  
77 other, even if the data points have low influence on the actual model calibration. To address this issue  
78 the modeller may choose to increase the robustness of the optimisation; however, these efforts will  
79 compound the computational demands of the case-deletion re-calibrations.

80 In regression applications Cook's distance can alternatively be calculated using "regression-theory"  
81 influence diagnostics (see Figure 1). Regression-theory influence diagnostics have a significantly  
82 reduced computational demand as they do not require case-deletion re-calibration and instead rely  
83 on assumptions about the type of regression model (linear/nonlinear) and residual error model  
84 (Gaussian, homoscedastic etc.). The reduced computational demand is achieved by combining the  
85 following two components for each observed data point: (1) the leverage, which describes the rate of

86 change of the predicted model output with respect to the corresponding observed output and can be  
87 used to assess the potential importance of individual observations [Wei *et al.*, 1998], and (2) the  
88 standardised residuals, which correspond to the raw residuals divided by the fitted standard deviation.  
89 By combining these two components to calculate Cook's distance, regression-theory influence  
90 diagnostics do not require additional re-calibrations and are therefore a more efficient alternative to  
91 the computationally demanding case-deletion influence diagnostics. There exist multiple alternative  
92 formulations of leverage, differing in the assumptions made about the fitted model and the  
93 probabilistic model of the residual errors. In circumstances where these assumptions are not violated  
94 regression-theory Cook's distance is equivalent to case-deletion Cook's distance.

95 Linear leverage is arguably the most widely used approach to approximate Cook's distance in  
96 regression problems [Fox and Weisberg, 2011], and is derived from standard linear regression theory  
97 and therefore inherits the assumptions of a linear model response (with respect to the model  
98 parameters) and Gaussian, homoscedastic and independent residual errors [Cook and Weisberg,  
99 1982]. When linear leverage is used in regression-theory Cook's distance (hereafter referred to as  
100 "linear Cook's distance") it also inherits these assumptions. This implies that linear Cook's distance  
101 may not be suitable for identifying the influential points in a hydrological modelling context as the  
102 hydrological model calibration violates the assumptions of linear regression, as a result of: 1) nonlinear  
103 model response [e.g. see discussion in Kavetski and Kuczera, 2007], and 2) heteroscedastic and non-  
104 Gaussian residual errors [e.g. see Schoups and Vrugt, 2010].

105 To address these limitations and expand the applicability of regression-theory influence diagnostics to  
106 more complex situations, St. Laurent and Cook [1992] proposed nonlinear leverage. Calculating Cook's  
107 distance by applying nonlinear leverage (hereafter referred to as "nonlinear Cook's distance") can take  
108 into account nonlinear model response, and is suitable for nonlinear models with Gaussian residuals.  
109 Wright *et al.* [2015] applied both linear and nonlinear Cook's distance in a hydrological modelling  
110 context and found that nonlinear Cook's distance provided higher performance than linear Cook's

111 distance, in terms of a higher correlation with the influential points identified using case-deletion  
112 influence diagnostics. The limitation of *Wright et al.* [2015] is that the hydrological models were  
113 calibrated using a standard least squares objective function, which is known to perform poorly in a  
114 hydrological modelling context when the residual errors are non-Gaussian and/or heteroscedastic [see  
115 *McInerney et al.*, 2017].

116 To overcome the limitations of the assumptions of linear and nonlinear leverage, generalised leverage  
117 was developed by *Wei et al.* [1998]. Generalised leverage makes no assumptions of linear model  
118 response, and can be applied to a broad range of objective functions, including those with  
119 heteroscedastic and/or non-Gaussian residual error assumptions. It has been used in numerous  
120 regression applications [e.g. *Leiva et al.*, 2014; *Lemonte and Bazán*, 2015; *Osorio*, 2016; *Rocha and*  
121 *Simas*, 2011]; however, it has not been applied in the context of hydrological or other environmental  
122 modelling applications. Furthermore, generalised leverage is typically used as a standalone diagnostic  
123 and has not previously been applied as an input to calculate Cook's distance (hereafter referred to as  
124 "generalised Cook's distance") to identify influential points. This research gap presents an opportunity  
125 to determine if generalised Cook's distance can be used as an efficient approach to approximate case-  
126 deletion Cook's distance in a computationally frugal manner.

127 Given the substantial computational advantages of regression-theory influence diagnostics over case-  
128 deletion influence diagnostics, they show significant promise for application in the field of hydrological  
129 and other environmental modelling applications. However, before regression-theory influence  
130 diagnostics can be applied, the validity of the assumptions of the formulations of leverage will first  
131 need to be experimentally tested in the context of hydrological case-studies. An important issue to be  
132 investigated is the hypothesis that generalised leverage can be used to approximate case-deletion  
133 Cook's distance as it has not previously been combined with standardised residuals to measure the  
134 proposed generalised Cook's distance. This study will assess the performance of the different  
135 approaches within the class of regression-theory influence diagnostics (i.e. linear Cook's distance,



136 nonlinear Cook's distance, and generalised Cook's distance) to reproduce case-deletion Cook's  
137 distance. The specific objectives of this study are to evaluate the ability of regression-theory influence  
138 diagnostics to identify influential points under the following modelling scenarios:

- 139 1. Linear and nonlinear regression models with either homoscedastic or heteroscedastic residual  
140 error;
- 141 2. A daily hydrological model including nonlinear model response and storage with  
142 heteroscedastic residual error; and
- 143 3. A stage-discharge rating curve model with Bayesian objective functions that include  
144 heteroscedastic residual error, data uncertainty and prior information.

145 For all three objectives, the regression-theory Cook's distance obtained using the linear, nonlinear and  
146 generalised leverage formulations will be compared to the case-deletion Cook's distance, in order to  
147 evaluate the extent to which the specific leverage formulation affects the performance of regression-  
148 theory influence diagnostics. The remainder of this paper is structured as follows. In Section 2 we  
149 describe the methodology, in Section 3 we introduce the three case studies selected to address the  
150 study objectives, and in Section 4 we apply the influence diagnostics to these case studies. In Section  
151 5 we discuss the advantages and disadvantages of case-deletion and regression-theory influence  
152 diagnostics, the suitability of applying generalised Cook's distance to a broader class of hydrological  
153 and environmental models, and the future need to understand the key drivers of influential data.

## 154 **2. Methodology**

155 Influence diagnostics identify data points that exert a disproportionate impact on calibrated  
156 parameters, performance and/or predictions. In this study we consider the following classes of Cook's  
157 distance influence diagnostics:

- 158 1. Case-deletion based Cook's distance, which measures the influence of a single point by  
 159 comparing model parameters, performance and/or predictions from calibration with and  
 160 without that data point; and
- 161 2. Regression-theory influence diagnostics, which measure influence by combining the  
 162 standardised residual and the leverage of each data point. We analyse and compare three  
 163 approaches to determining the leverage, which produce three estimates of Cook's distance:
- 164 i. Linear Cook's distance, which uses linear leverage,
  - 165 ii. Nonlinear Cook's distance, which uses nonlinear leverage, and
  - 166 iii. Generalised Cook's distance, which uses generalised leverage.

167 In this section we first introduce the general modelling framework, and then define the influence  
 168 diagnostics, leverage, and the objective functions used in this study. We finish by describing the  
 169 metrics that we will use to evaluate the performance of the regression-theory influence diagnostics.

## 170 **2.1. General model framework**

171 We define the general model response as:

$$172 \quad \mathbf{y} = f(\boldsymbol{\alpha}; \mathbf{X}) + \boldsymbol{\varepsilon} \quad (1)$$

173 where  $\mathbf{y} = (y_1, y_2, \dots, y_n)$  is a vector of  $n$  observed responses,  $f(\cdot)$  is the model structure,  
 174  $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_{m_\alpha})$  is a vector of  $m_\alpha$  model parameters,  $\mathbf{X}$  is an  $n \times k$  matrix of  $k$  observed inputs  
 175 (e.g., precipitation, potential evapotranspiration (PET)), and  $\boldsymbol{\varepsilon}$  is a vector of  $n$  residual errors.  
 176 Residuals are further assumed to be realisations from a given probability distribution, with parameters  
 177  $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_{m_\beta})$  (e.g. a centred Gaussian distribution with unknown standard deviation). Thus,  
 178 the entire set of  $m$  parameters to be calibrated are  $\boldsymbol{\theta} = \{\boldsymbol{\alpha}, \boldsymbol{\beta}\}$  which includes both the model  
 179 parameters  $\boldsymbol{\alpha}$  and the residual error model parameters  $\boldsymbol{\beta}$ .

### 180           **2.1.1. Objective functions**

181   In order to apply leverage to a broad class of objective functions used in hydrological modelling we  
182   consider the general form of the objective function, as suggested by *Wei et al.* [1998]:

$$183 \quad \Phi(\boldsymbol{\theta}; \mathbf{y}, \mathbf{X}) = \sum_{i=1}^n \rho_i(f_i(\boldsymbol{\alpha}; \mathbf{X}), \boldsymbol{\beta}; y_i) \quad (2)$$

184   where  $\rho_i(\cdot)$  is a function that describes the contribution of the  $i^{\text{th}}$  data point to the objective  
185   function,  $f_i(\boldsymbol{\alpha}; \mathbf{X})$  is the  $i^{\text{th}}$  model prediction,  $\Phi(\boldsymbol{\theta}; \mathbf{y}, \mathbf{X})$  and  $f(\boldsymbol{\alpha}; \mathbf{X})$  are assumed to be twice  
186   differentiable with respect to  $\boldsymbol{\theta}$  and  $\mathbf{y}$ . We will denote  $\hat{\boldsymbol{\theta}}$  as the model parameters that maximise  $\Phi$   
187   in equation (2), and  $\hat{\mathbf{y}}$  as the predicted response associated with  $\hat{\boldsymbol{\theta}}$ , i.e.  $\hat{\mathbf{y}} = f(\hat{\boldsymbol{\alpha}}; \mathbf{X})$ .

188   The generalised form in equation (2) can be adapted to a number of well-known objective functions  
189   in hydrological modelling as outlined in Section 2.4.

### 190           **2.1.2. Standardised residuals**

191   The standardised residuals,  $\mathbf{v}$ , which are required to estimate the regression-theory influence  
192   diagnostics introduced in Section 2.2.2, are obtained by dividing the raw residuals  $\boldsymbol{\varepsilon} = \mathbf{y} - \hat{\mathbf{y}}$  by their  
193   calibrated standard deviations,  $\boldsymbol{\sigma}$ :

$$194 \quad \mathbf{v} = \frac{\boldsymbol{\varepsilon}}{\boldsymbol{\sigma}} \quad (3)$$

195   The vector  $\boldsymbol{\sigma}$  is determined based on the assumed residual error model and the resultant objective  
196   function (see Section 2.4 for further details).

## 197           **2.2. Influence diagnostics**

198 This section provides a detailed description of the different influence diagnostics used in this study  
 199 (see Figure 1 for an overview). Firstly, we present the case-deletion “class” of influence diagnostics  
 200 and outline the approach used to calculate case-deletion Cook’s distance. Secondly, we present the  
 201 regression-theory “class” of influence diagnostics and outline the approaches used to calculate  
 202 regression-theory Cook’s distance using three formulations of leverage (i.e. linear, nonlinear and  
 203 generalised leverage) to produce linear, nonlinear and generalised Cook’s distance.

### 204 **2.2.1. Case-deletion influence diagnostics**

205 Case-deletion influence diagnostics describe the influence of masking a data point in model calibration  
 206 and assessing the change to the model predictions, parameters and/or objective function value.  
 207 Cook’s distance can be measured exactly using case-deletion [see *Cook and Weisberg, 1982*]; note  
 208 that in the statistical literature this case-deletion Cook’s distance is sometimes referred to as  
 209 “generalised Cook’s distance” [e.g. *Das, 2008*]. Case-deletion based Cook’s distance measures  
 210 influence by comparing model predictions  $\mathbf{y}$  based on using all of the calibration data and model  
 211 predictions  $\mathbf{y}^{(-i)}$  with the  $i^{\text{th}}$  point masked from the calibration data. For a given data point, case-  
 212 deletion based Cook’s distance is calculated by:

$$213 \quad CD_i = \sum_{j=1}^n \frac{(\hat{y}_j - \hat{y}_j^{(-i)})^2}{m \times \hat{\sigma}_j^2} \quad (4)$$

214 where  $\sigma_j$  is the calibrated standard deviation for the  $j^{\text{th}}$  data point, estimated from using all  
 215 calibration data (i.e.  $\mathbf{y}$ ).

### 216 **2.2.2. Regression-theory influence diagnostics**

217 Regression-theory influence diagnostics avoid the computational burden of case-deletion re-  
 218 calibration by making assumptions about the type of response model (linear/nonlinear) and residual

219 error model (Gaussian, homoscedastic etc.). Regression-theory Cook's distance is calculated by  
 220 combining the standardised residual of the  $i^{\text{th}}$  point ( $v_i$ ) with the leverage of  $i^{\text{th}}$  observation on the  
 221  $i^{\text{th}}$  prediction ( $L_{ii}$ ) to give [Cook and Weisberg, 1982; Fox and Weisberg, 2011]:

$$222 \quad CD_i = \frac{v_i^2}{m} \frac{L_{ii}}{(1-L_{ii})^2} \quad (5)$$

223 The approach used to determine the three different forms of Cook's distance (i.e. linear, nonlinear  
 224 and generalised Cook's distance; Figure 1) is based on the corresponding forms of leverage (i.e. linear,  
 225 nonlinear, and generalised leverage). In the next section, we provide a general definition of leverage  
 226 followed by the three specific formations of leverage that are used to calculate regression-theory  
 227 Cook's distance.

### 228 **2.3. Leverage**

229 Leverage generally can be defined as the rate of the change of the  $i^{\text{th}}$  predicted value,  $y_i$ , with  
 230 respect to another  $j^{\text{th}}$  observed value,  $y_j$  [Cook and Weisberg, 1982; Hoaglin and Welsch, 1978; St.  
 231 Laurent and Cook, 1992; Wei et al., 1998]:

$$232 \quad L_{ij} = \partial \hat{y}_i / \partial y_j \quad (6)$$

233 or in matrix notation:

$$234 \quad \mathbf{L} = \frac{\partial \mathbf{y}}{\partial \mathbf{y}^T} \quad (7)$$

235 where  $\mathbf{L}$  is an  $n \times n$  matrix. The diagonal elements  $L_{ii}$  most directly reflect the impact of  $y_i$  on the  
 236 model fit [Cook and Weisberg, 1982; Hoaglin and Welsch, 1978; St. Laurent and Cook, 1992], and are  
 237 used for calculating regression-theory Cook's distance (Section 2.2.2).

238 **2.3.1. Linear leverage**

239 Linear leverage inherits the assumptions of standard linear regression theory; i.e. that the model  
 240 response (with respect to the parameters) is linear and that residual errors are Gaussian,  
 241 homoscedastic and independent. Under the assumptions of linear regression the general form of  
 242 leverage in equation (7) can be expressed as  $\mathbf{L}$  [Fox and Weisberg, 2011]:

243 
$$\mathbf{L} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \quad (8)$$

244 As linear leverage depends solely on the observed input  $\mathbf{X}$ , it can be calculated without model  
 245 calibration using linear algebra. In a linear regression model with standard least squares (SLS) residual  
 246 errors, regression-theory Cook's distance is equivalent to case-deletion Cook's distance [see Cook,  
 247 1977].

248 **2.3.2. Nonlinear leverage**

249 Nonlinear leverage does not assume a linear model response but retains the assumption that residual  
 250 errors are Gaussian, homoscedastic and independent. Nonlinear leverage is dependent on the local  
 251 sensitivity of the model predictions to small perturbations in model parameters [St. Laurent and Cook,  
 252 1992]. Nonlinear leverage is calculated after model calibration, and under the assumptions of  
 253 nonlinear regression the general form of leverage in equation (7) can be expressed as  $\mathbf{L}(\boldsymbol{\alpha})$  [St.  
 254 Laurent and Cook, 1992; 1993; Wei et al., 1998; Wright et al., 2015]:

255 
$$\mathbf{L}(\boldsymbol{\alpha}) = \frac{\partial f(\boldsymbol{\alpha}; \mathbf{X})}{\partial \boldsymbol{\alpha}} \left( \left( \frac{\partial f(\boldsymbol{\alpha}; \mathbf{X})}{\partial \boldsymbol{\alpha}} \right)^T \frac{\partial f(\boldsymbol{\alpha}; \mathbf{X})}{\partial \boldsymbol{\alpha}} - \sum_{i=1}^n \left( (y_i - \hat{y}_i) \frac{\partial^2 f_i(\boldsymbol{\alpha}; \mathbf{X})}{\partial \boldsymbol{\alpha}^2} \right) \right)^{-1} \left( \frac{\partial f(\boldsymbol{\alpha}; \mathbf{X})}{\partial \boldsymbol{\alpha}} \right)^T \quad (9)$$

256 where  $\frac{\partial f(\boldsymbol{\alpha}; \mathbf{X})}{\partial \boldsymbol{\alpha}}$  is the  $n \times m_\alpha$  Jacobian matrix with  $i^{\text{th}}$  row  $\frac{\partial f_i(\boldsymbol{\alpha}; \mathbf{X})}{\partial \boldsymbol{\alpha}}$ , and  $\frac{\partial^2 f_i(\boldsymbol{\alpha}; \mathbf{X})}{\partial \boldsymbol{\alpha}^2}$  is the  
 257  $m_\alpha \times m_\alpha$  Hessian matrix associated with the  $i^{\text{th}}$  data point. Analytical derivatives are typically not

258 available for hydrological models, and therefore we obtain estimates of the derivatives from central-  
 259 difference numerical approximation [Nocedal and Wright, 2006]. When applied to a linear regression  
 260 model with SLS residual errors, the nonlinear leverage simplifies to linear leverage, as shown in *Wei*  
 261 *et al.* [1998].

### 262 **2.3.3. Generalised leverage**

263 Generalised leverage makes no assumptions of linear model response, and can be applied to a general  
 264 class of regression models and a broad range of objective functions, including those with  
 265 heteroscedastic and/or non-Gaussian residual error assumptions. Generalised leverage is calculated  
 266 after model calibration and takes into account the curvature of the objective function about the whole  
 267 set of calibrated parameters  $\boldsymbol{\theta}$ . In this case the general form of leverage in equation (7) can be  
 268 expressed as  $\mathbf{L}(\boldsymbol{\theta})$  [Wei *et al.*, 1998]:

$$269 \quad \mathbf{L}(\boldsymbol{\theta}) = \frac{\partial f(\boldsymbol{\alpha}; \mathbf{X})}{\partial \boldsymbol{\theta}} \left( -\frac{\partial^2 \Phi(\boldsymbol{\theta}; \mathbf{y}, \mathbf{X})}{\partial \boldsymbol{\theta}^2} \right)^{-1} \frac{\partial^2 \Phi(\boldsymbol{\theta}; \mathbf{y}, \mathbf{X})}{\partial \boldsymbol{\theta} \partial \mathbf{y}^T} \quad (10)$$

270 where  $\frac{\partial f(\boldsymbol{\alpha}; \mathbf{X})}{\partial \boldsymbol{\theta}}$  is the  $n \times m$  Jacobian matrix with  $i^{\text{th}}$  row  $\frac{\partial f_i(\boldsymbol{\alpha}; \mathbf{X})}{\partial \boldsymbol{\theta}}$  (note that  $\frac{\partial f_i(\boldsymbol{\alpha}; \mathbf{X})}{\partial \boldsymbol{\beta}} = 0$ ),

271  $\frac{\partial^2 \Phi(\boldsymbol{\theta}; \mathbf{y}, \mathbf{X})}{\partial \boldsymbol{\theta}^2}$  is a  $m \times m$  Hessian matrix and  $\frac{\partial^2 \Phi(\boldsymbol{\theta}; \mathbf{y}, \mathbf{X})}{\partial \boldsymbol{\theta} \partial \mathbf{y}^T}$  is a  $m \times n$  matrix. Generalised leverage can

272 be applied to any objective function that takes the general form in equation (2). Generalised leverage  
 273 simplifies to nonlinear leverage in the case of a nonlinear regression model and SLS residual errors, as  
 274 shown in *Wei et al.* [1998].

### 275 **2.4. Objective functions used in this study**

276 This section introduces the range of different objective functions that will be used in the case studies  
 277 to evaluate the performance of the differing implementations of regression-theory Cook's distance.

### 278 **2.4.1. Standard least squares**

279 Assuming independent and identically distributed (i.i.d.) Gaussian residual errors, the following log  
280 likelihood can be used as an objective function:

$$281 \quad \Phi(\boldsymbol{\theta}; \mathbf{y}, \mathbf{X}) = \sum_{i=1}^n \log(p_N(y_i - f_i(\boldsymbol{\alpha}; \mathbf{X}) | 0, \sigma^2)) \quad (11)$$

282 where  $p_N(x | \mu, \sigma^2)$  is the Gaussian probability density at  $x$  assuming constant mean  $\mu$  and  
283 variance  $\sigma^2$ . As the standard deviation  $\sigma$  is unknown it will be estimated, and therefore we have  
284  $\beta = \{\sigma\}$ . Note that the Nash-Sutcliffe efficiency [Nash and Sutcliffe, 1970] objective function that is  
285 commonly applied in hydrological calibration corresponds to the assumptions of constant-variance  
286 and Gaussian residual errors of the standard least squares (SLS) objective function. Note that (11) is a  
287 particular case of the general objective function in equation (2).

### 288 **2.4.2. Weighted least squares**

289 Residual errors in hydrological applications are generally heteroscedastic [see Schoups and Vrugt,  
290 2010; Sorooshian and Dracup, 1980] and to account for this non-constant variance we apply a  
291 weighted least squares (WLS) objective function. Due to this heteroscedasticity in hydrological  
292 residual errors it is common to replace the constant standard deviation  $\sigma$  in equation (11) with a  
293 standard deviation  $\sigma$  that varies in time, so that the non-constant variance acts as a “weight” for each  
294 residual [e.g. McInerney et al., 2017; Thyer et al., 2009]. A common covariate for modelling  
295 heteroscedasticity in streamflow errors is the predicted streamflow itself [e.g. Schoups and Vrugt,  
296 2010; Thyer et al., 2009]. Following Evin et al. [2014] we consider the standard deviation of residuals  
297 to be a linear function of simulated streamflow, such that:

$$298 \quad \sigma = \beta_1 \mathbf{y} + \beta_2 \quad (12)$$

299 The objective function becomes:



300 
$$\Phi(\boldsymbol{\theta}; \mathbf{y}, \mathbf{X}) = \sum_{i=1}^n \log(p_N(y_i - f_i(\boldsymbol{\alpha}; \mathbf{X}) | 0, \sigma_i^2)) \quad (13)$$

301 As the parameters describing the non-constant standard deviation (i.e.  $\boldsymbol{\beta} = \{\beta_1, \beta_2\}$ ) are unknown  
 302 they will need to be estimated. Note that (12) is a particular case of the general objective function in  
 303 equation (2).

304 **2.4.3. Weighted least squares with data uncertainty**

305 In circumstances when independent estimates of data errors are available we may wish to distinguish  
 306 between heteroscedasticity in hydrological residual errors and uncertainty of observed responses. To  
 307 implement the WLS method with discharge uncertainty in the WLS objective function (13) we assume  
 308 that the total errors can be decomposed as the sum of two independent error terms: the “structural  
 309 errors” that can be described using the WLS standard deviation  $\boldsymbol{\sigma}_r = \beta_1 \mathbf{y} + \beta_2$  and the “measurement  
 310 errors” described using known standard deviations  $\boldsymbol{\sigma}_y$ . The latter standard deviations may be derived  
 311 from an uncertainty analysis of measured responses, which can be performed before and  
 312 independently from the model calibration. The standard deviation of the total error, combining  
 313 structural and measurement errors, is therefore equal to  $\boldsymbol{\sigma} = \sqrt{\boldsymbol{\sigma}_r^2 + \boldsymbol{\sigma}_y^2}$ . Hence the  $\sigma_i$  in equation  
 314 (13) becomes:

315 
$$\sigma_i = \sqrt{(\beta_1 y_i + \beta_2)^2 + \sigma_{y,i}^2} \quad (14)$$

316 where  $\sigma_{y,i}$  is the standard deviation of the measurement errors at time step  $i$ .

317 **2.4.4. Weighted least squares with priors**

318 In circumstances when prior information about parameter values is available based on previous  
 319 studies and/or from analysis of physical characteristics that govern the relation between inputs  $\mathbf{X}$

320 and outputs  $\mathbf{y}$  we can use an objective function that combines WLS likelihood with priors. Bayes'  
 321 equation yields the posterior probability distribution of the hydrological and residual error model  
 322 parameters as follows:

$$323 \quad \underbrace{p(\boldsymbol{\theta}|\mathbf{X},\mathbf{y})}_{\text{posterior}} \propto \underbrace{p(\mathbf{y}|\boldsymbol{\theta},\mathbf{X})}_{\text{likelihood}} \underbrace{p(\boldsymbol{\theta})}_{\text{prior}} \quad (15)$$

324 where  $p(\boldsymbol{\theta}|\mathbf{X},\mathbf{y})$  is the posterior probability of parameter  $\boldsymbol{\theta}$  given  $\mathbf{X}$  and  $\mathbf{y}$ ,  $p(\boldsymbol{\theta})$  is the joint prior  
 325 probability density of hydrological and residual error model parameters, and  $p(\mathbf{y}|\boldsymbol{\theta},\mathbf{X})$  is the  
 326 likelihood of  $\mathbf{y}$  given  $\boldsymbol{\theta}$  and  $\mathbf{X}$ . Taking the logarithm of equation (15) we obtain:

$$327 \quad \log(p(\boldsymbol{\theta}|\mathbf{X},\mathbf{y})) = \log(p(\mathbf{y}|\boldsymbol{\theta},\mathbf{X})) + \log(p(\boldsymbol{\theta})) + c \quad (16)$$

328 where  $c$  is a constant. Assuming independence between residuals we can formulate the objective  
 329 function as:

$$\begin{aligned} \Phi(\boldsymbol{\theta};\mathbf{y},\mathbf{X}) &= \log(p(\mathbf{y}|\boldsymbol{\theta},\mathbf{X})) + \log(p(\boldsymbol{\theta})) \\ 330 \quad &= \sum_{i=1}^n \log(p(y_i|\boldsymbol{\theta},\mathbf{X})) + \log(p(\boldsymbol{\theta})) \\ &= \sum_{i=1}^n \left( \log(p(y_i|\boldsymbol{\theta},\mathbf{X})) + \frac{1}{n} \log(p(\boldsymbol{\theta})) \right) \end{aligned} \quad (17)$$

331 Assuming the residual errors are heteroscedastic with  $\sigma$  given by equation (12) and independent  
 332 priors, we obtain the following objective function:

$$333 \quad \Phi(\boldsymbol{\theta};\mathbf{y},\mathbf{X}) = \sum_{i=1}^n \left\{ \log(p_N(y_i - f_i(\boldsymbol{\alpha};\mathbf{X}) | 0, \sigma_i^2)) + \frac{1}{n} \sum_{j=1}^p \log(p(\theta_j)) \right\} \quad (18)$$

334 where the contributions to the objective function from the priors are split evenly across the  $n$  points  
 335 in the calibration data.

#### 336 **2.4.5. Weighted least squares with data uncertainty and priors**

337 In circumstances when both independent estimates of data errors and prior information about  
338 parameter values are available we can use weighted least squares with data uncertainty and priors.  
339 Similar to Section 2.4.3, data uncertainty can readily be included in the weighted least squares with  
340 priors objective function (18) by simply using  $\sigma = \sqrt{\sigma_r^2 + \sigma_y^2}$ , where  $\sigma_r = \beta_1 y + \beta_2$ , and  $\sigma_y$  are  
341 known values representing the measurement uncertainty in observed responses.

## 342 **2.5. Performance metrics**

343 As case-deletion Cook's distance provides a measure of influence with no assumptions regarding the  
344 type of model (linear/nonlinear) or the complexity of the residual error model (Gaussian,  
345 heteroscedastic, etc.) we use it as a baseline to compare the three formulations of regression-theory  
346 influence diagnostics: linear Cook's distance, nonlinear Cook's distance and generalised Cook's  
347 distance. We use two metrics to assess the performance of regression-theory influence diagnostics  
348 with respect to case-deletion based Cook's distance. These metrics are evaluated on 1) the whole set  
349 of influential data points, to show the general ability of regression-theory influence diagnostics to  
350 approximate case-deletion Cook's distance; and 2) a subset comprising the 10 most influential data  
351 points identified by case-deletion Cook's distance, to highlight the performance with respect to the  
352 points that are most influential to calibration. The metrics are:

- 353 1. Spearman correlation (Sp. and Sp.<sub>10</sub>), which provides a measure of the performance of the  
354 regression-theory influence diagnostics to correctly rank the most influential data points.
- 355 2. Coefficient of determination ( $r^2$  and  $r^2_{10}$ ), which provides a measure of the proportion of the  
356 variance in the case-deletion based variable that is accounted for by the regression-theory  
357 variable.

358 The selected performance metrics allow for a thorough comparison of the regression-theory influence  
359 diagnostics as approximations of the case-deletion Cook's distance.

## 360 **3. Case studies**

361 The research objectives of this paper are to evaluate the ability of regression-theory influence  
362 diagnostics to identify influential points under the following modelling scenarios:

- 363 1. Linear and nonlinear regression models with either homoscedastic or heteroscedastic residual  
364 error;
- 365 2. A daily hydrological model including nonlinear model response and storage with  
366 heteroscedastic residual error; and
- 367 3. A stage-discharge rating curve model with Bayesian objective functions that include  
368 heteroscedastic residual error, data uncertainty and prior information.

369 In order to address these objectives we apply case-deletion and regression-theory influence  
370 diagnostics to ten different case studies, organised in three distinct case study sets (Table 1). To  
371 address the first research objective the first case study set consists of four synthetic regression  
372 models,  $A_{1-4}$ , are selected to test the performance with linear/nonlinear regression models and  
373 homoscedastic/heteroscedastic residual error models. The second research objective is addressed by  
374 case study set 2, which tests the performance with daily hydrological models,  $B_{1-2}$ , with nonlinear  
375 hydrological response, model storage, and heteroscedastic residual errors. Finally, the third objective  
376 is addressed by case study set 3, which tests the performance with four different rating curve models,  
377  $C_{1-4}$ , with and without data uncertainty and with and without prior knowledge specified using a  
378 Bayesian inference approach.

379 In all cases the objective functions are optimised using the Shuffled Complex Evolution (SCE) search  
380 algorithm [Duan *et al.*, 1992; Duan *et al.*, 1994] followed by a Nelder-Mead gradient search from the  
381 SCE optimised parameter set to machine precision to ensure convergence to the optima.

### 382 **3.1. Case study set 1: Synthetic regression models with linear/nonlinear** 383 **response and homoscedastic/heteroscedastic residual errors**

384 The first case study set uses synthetic regression models that range in complexity from a simple linear  
385 model response with homoscedastic residual errors to a nonlinear power model response with  
386 heteroscedastic residual errors. The regression models with synthetic data ( $A_{1-4}$ ; Table 1) are selected  
387 to highlight the role of model structure and residual error model on the influence results:  $A_1$  has a  
388 linear model response with a standard least squares (SLS) residual error model;  $A_2$  also has a linear  
389 model response but with a weighted least squares (WLS) residual error model; and both  $A_3$  and  $A_4$   
390 have a nonlinear model response with SLS and WLS residual error, respectively.

### 391 **3.2. Case study set 2: Daily hydrological model with synthetic and observed** 392 **streamflow and heteroscedastic residual errors**

393 The next case study set tests the performance of the regression-theory influence diagnostics in a  
394 typical hydrological modelling calibration context. We apply a daily hydrological model that includes  
395 nonlinear model response and storage (meaning that inputs at a given time-step can affect outputs  
396 many time-steps into the future) and heteroscedastic residual errors. The daily lumped hydrological  
397 model GR4J [Perrin *et al.*, 2003] was selected based upon its popularity [e.g. Andréassian *et al.*, 2014;  
398 Evin *et al.*, 2014; Le Moine *et al.*, 2007; Lebecherel *et al.*, 2016; Wright *et al.*, 2015] and parsimonious  
399 model structure. This allows for computational efficiency in the case-deletion model runs required to  
400 calculate case-deletion Cook's distance. The GR4J hydrological model has model parameters  
401  $\alpha = \{\alpha_1, \alpha_2, \alpha_3, \alpha_4\}$ , where  $\alpha_1$  is the maximum capacity of the production store,  $\alpha_2$  is the  
402 groundwater exchange coefficient,  $\alpha_3$  is the maximum capacity of the routing store, and  $\alpha_4$  is the  
403 time base of unit hydrograph.

404 We apply the GR4J hydrological model to the French Broad River catchment in North Carolina, USA.  
405 The French Broad River has a catchment area of 2448 km<sup>2</sup>, annual precipitation of 1413 mm and  
406 annual streamflow of 800 mm, leading to a runoff coefficient of 0.57.

407 We explore two alternative modelling scenarios  $B_{1-2}$  (Table 1) that correspond to synthetic streamflow  
 408 data and real observed streamflow data, respectively. We use three years of calibration data, from  
 409 1974 to 1976. The first model  $B_1$  uses the observed rainfall and PET from the French Broad River but  
 410 has synthetic streamflow data. This synthetic streamflow data is obtained by first using real  
 411 streamflow data to fit the GR4J parameters, then using the fitted parameters to generate a predicted  
 412 streamflow time series, and finally adding residual errors to the predicted time series based on the  
 413 WLS error model. The second hydrological model  $B_2$  also uses observed rainfall and PET from the  
 414 French Broad River catchment, but is calibrated to the real observed streamflow data. Note that while  
 415 there are two inputs for GR4J (i.e. rainfall and PET), here we consider only the importance of rainfall  
 416 data (i.e. don't include PET in  $\mathbf{X}$ ) when calculating leverage, because typically hydrological model  
 417 response are more sensitive to errors in rainfall, rather than errors in PET [Oudin *et al.*, 2006].

### 418 **3.3. Case study set 3: Rating curve model incorporating heteroscedastic residual** 419 **errors, data uncertainty and parameter priors**

420 The final case study set uses a rating curve model, with increasing complexity in the objective function  
 421 that investigates the impact of data uncertainty and incorporating parameter priors using a Bayesian  
 422 approach. We apply a piecewise stage-discharge rating curve model to the Ardèche River at Sauze,  
 423 France. The Ardèche River has a catchment area of 2240 km<sup>2</sup> with a mean annual discharge of 63 m<sup>3</sup>/s.  
 424 We use the reduced subset of 38 stage-discharge gaugings applied in *Le Coz et al.* [2014]. The flow at  
 425 the hydrometric station is controlled by a rectangular sill at low flows, and a rectangular channel at  
 426 high flows, leading to a two-part rating curve model with the following stage-discharge relationship:

$$427 \quad f(\mathbf{a}, X_i) = \begin{cases} a_1 (X_i - b_1)^{c_1}, & \text{for } X_i < k \\ a_2 (X_i - b_2)^{c_2}, & \text{for } X_i \geq k \end{cases} \quad (19)$$

428 Here  $\mathbf{X}$  is stage and  $\mathbf{a} = \{a_1, b_1, c_1, k, a_2, c_2\}$  are the rating curve model parameters similar to *Le Coz*  
 429 *et al.* [2014]. As the rating curve is continuous at the knot ( $k$ ), the parameter  $b_2$  is computed from

430 the other calibrated parameter values by solving the continuity condition  $a_1(k - b_1)^{c_1} = a_2(k - b_2)^{c_2}$   
431 , yielding  $b_2 = k - \left( (a_1 / a_2) (k - b_1)^{c_1} \right)^{1/c_2}$ . Petersen-Øverleir [2004] suggest a heteroscedastic  
432 residual error model to take into account the heteroscedasticity of most rating curve errors, and so  
433 we use the WLS objective function described in Section 2.4.2. We apply the following four calibration  
434 schemes across  $C_{1-4}$ : 1) baseline rating curve calibration with WLS in  $C_1$ ; 2) rating curve calibration  
435 with discharge uncertainty in  $C_2$ ; 3) rating curve with priors in  $C_3$ ; and 4) rating curve calibration with  
436 discharge uncertainty and priors in  $C_4$ .

437 We follow *Le Coz et al.* [2014] who provide gauging uncertainties for the discharge data at Sauze and  
438 also a framework for Bayesian inference. In  $C_3$  and  $C_4$  we use the priors from *Le Coz et al.* [2014] for  
439 the model parameters that are summarised in Table 2. Perusal of Table 2 shows that the prior standard  
440 deviation is smallest for the exponent parameters ( $c_1$  and  $c_2$  in equation (19)), compared with the  
441 scaling parameters,  $a_1$  and  $a_2$ , and the offset parameters,  $b_1$  and  $b_2$ . Hence the priors are more  
442 informative for these exponent values because they only depend on the type of hydraulic control  
443 (here, rectangular sill and rectangular channel). In the case of the residual error model parameters  $\beta$   
444 there is no prior knowledge and so an uninformative uniform distribution is applied.

#### 445 **4. Assessing the ability of regression-theory Cook's distance to reproduce case-** 446 **deletion Cook's distance**

447 We apply case-deletion and regression-theory influence diagnostics with linear, nonlinear and  
448 generalised Cook's distance to the three case studies in Sections 4.1-4.3. In Section 4.4 we summarise  
449 the performance of the regression-theory influence diagnostics across the case studies, and we finish  
450 in Section 4.5 with an analysis of the computation times of both the regression-theory and case-  
451 deletion based influence diagnostics.

#### 4.1. Case study set 1: Synthetic regression models with linear/nonlinear response and homoscedastic/heteroscedastic residual errors

In this section we evaluate the performance of regression-theory Cook's distance based on the three formulations of generalised leverage, using synthetic regression case studies with varying degrees of nonlinear model response and heteroscedastic residual errors ( $A_{1-4}$ ; Table 1). The synthetically generated "observed" data and fitted models are presented in Figure 2 (row 1) for the four cases. The models are correctly specified, and fit the data well in all cases. This is evidenced by the standardised residuals being independent and normally distributed, with zero mean and unit standard deviation (Figure 2, row 2).

Similarities and differences between the three leverage formulations are shown in Figure 2 (row 3). Linear leverage is smooth and parabolic in all four cases ( $A_{1-4}$ ), with a minima at the mean of  $\mathbf{X}$  ( $\sim 100$ ). This highlights that linear leverage only depends on input  $\mathbf{X}$  (which is identical in all four cases), and therefore does not vary with the case study. Nonlinear leverage is the same as linear leverage for linear response models  $A_1$  and  $A_2$ , but differs for nonlinear response models  $A_3$  and  $A_4$ . In those cases, the nonlinear model response results in higher leverage for larger values of  $\mathbf{X}$ , with a slight increase in the midrange of  $\mathbf{X}$  for  $A_3$ , and with leverage varying smoothly as a function of  $\mathbf{X}$ . Interestingly, the nonlinear leverage for case  $A_3$  is different to the nonlinear leverage for  $A_4$ . This is due to slightly different calibrated parameter values  $\hat{\alpha}$  for the nonlinear model in  $A_3$  compared with  $A_4$ ; if these calibrated parameter values were identical, the nonlinear leverage in equation (9) would be the same, since it is a function of input data  $\mathbf{X}$ , model response  $f()$ , and optimal model parameters  $\hat{\alpha}$ . This highlights the sensitivity of nonlinear leverage to influential data points, despite the observations  $\mathbf{y}$  not appearing explicitly in equation (9). Finally, generalised leverage is the same as nonlinear leverage for cases  $A_1$  and  $A_3$ , when residuals are homoscedastic. However, when heteroscedasticity in residuals is introduced into the "observations" and likelihood functions (cases  $A_2$  and  $A_4$ ), we see there are two major differences. The first difference is that generalised leverage becomes larger than nonlinear



477 leverage for small values of  $\mathbf{X}$ . This is because generalised leverage accounts for the higher weights  
478 (i.e. smaller standard deviations) placed on low values of  $\mathbf{Y}$  in the WLS likelihood function (which  
479 correspond to small values of  $\mathbf{X}$ ), while nonlinear leverage applies the same weight to all values of  
480  $\mathbf{Y}$ . The second differences is that unlike linear and nonlinear leverage, generalised leverage does not  
481 vary smoothly as a function of  $\mathbf{X}$ . This is because for a given point  $i$ , the generalised leverage in  
482 equation (10) depends on the observation at that point  $y_i$ , and the observations  $\mathbf{y}$  do not vary  
483 smoothly with  $\mathbf{X}$ .

484 The magnitude of the case-deletion Cook's distance is presented in Figure 2 (row 4) as grey bubbles,  
485 and compared to the regression-theory Cook's distance (which combines the leverage and  
486 standardised residuals, equation (5)) in Figure 2 (row 5) as a function of  $\mathbf{X}$ . The differences between  
487 case-deletion Cook's distance and the three regression-theory Cook's distances are also quantified in  
488 Figure 3. The three regression-theory Cook's distances are identical for case  $A_1$ , as a result of identical  
489 leverages. The errors between the regression-theory Cook's distance and case-deletion Cook's  
490 distance are small (green, purple and orange bubbles are all similarly small in Figure 2, column 1, row  
491 5) and the correlations are high (as evidence by  $r^2$  values and Spearman correlations of 1.00 when  
492 calculated over all data and the top 10 most influential points in Figure 3, column 1).

493 When heteroscedastic residual errors are introduced (case  $A_2$ ), generalised Cook's distance becomes  
494 the most accurate approximation (green bubbles show lower errors than purple bubbles in Figure 2,  
495 column 2, row 5), with linear and nonlinear Cook's distance being the same (purple bubbles overlay  
496 orange bubbles). For linear and nonlinear Cook's distance, performance is worst at the extremes of  
497  $\mathbf{X}$ , and particularly the lower values of  $\mathbf{X}$  as they do not account for residual heteroscedasticity.  
498 The increased accuracy of using generalised Cook's distances is seen in the top 10% of influential  
499 points (Figure 3, column 2, row 2) where—relative to the other leverage formulations—the Spearman  
500 correlation increases from 0.65 to 0.96, and the  $r^2$  increases from 0.28 to 0.98.

501 The nonlinear response with homoscedastic residual errors (case  $A_3$ ) results show identical  
502 performance for the nonlinear and generalised Cook's distances, which are typically more accurate  
503 than linear Cook's distance (green and purple bubbles have lower errors than orange bubbles in Figure  
504 2, column 3, row 5). Linear Cook's distance performs particularly poorly for high values of  $\mathbf{X}$ , as  
505 anticipated based on the leverage results. The largest improvement is obtained by using nonlinear and  
506 generalised Cook's distances is seen in the top 10% of influential points (Figure 3, column 3, row 2)  
507 where the Spearman correlation increases from 0.75 to 1.00, and the  $r^2$  increases from 0.50 to 1.00.

508 Finally, the nonlinear model response with heteroscedastic residual errors (case  $A_4$ ) results show that  
509 the generalised Cook's distance is the most accurate of the regression-theory Cook's distances (green  
510 bubbles show the lowest error in Figure 2, column 3, row 5). Both Spearman correlation and  $r^2$  values  
511 are close to unity in all cases except for the Spearman correlation value for the largest 10% of  
512 influential points (Sp. = 0.79), due to a difference in a single point - the largest Cook's distance value.

513 The ranking of the performance linear and nonlinear Cook's distance for this case appears to depend  
514 on  $\mathbf{X}$  and the accuracy matrix used (abs. errors, correlation or Spearman rank on all or top 100 data  
515 points). Neither of these leverage approaches, produce the consistent accuracy of generalised Cook's  
516 distance.

517 Overall, the results indicate that for the four synthetic regression model case studies considered,  
518 generalised Cook's distance provides a very close approximation of case-deletion Cook's distance, and  
519 represents a significant improvement in identifying the influential points compared to the other  
520 regression-theory influence diagnostics linear Cook's distance and nonlinear Cook's distance.

#### 521 **4.2. Case study set 2: Daily hydrological model with synthetic and observed** 522 **streamflow and heteroscedastic residual errors**

523 We now evaluate the performance of regression-theory influence diagnostics in a typical hydrological  
524 modelling context where the model has nonlinear response, storage and heteroscedastic errors, with  
525 both synthetic and real observed catchment data (models  $B_1$  and  $B_2$ , respectively; see Table 1).

526 Observed and predicted streamflow is shown in the first row in Figure 4 for three representative time  
527 periods. For case  $B_1$ , when synthetic streamflow data is used for “observations”, the hydrological  
528 model provides a good fit to the observations for both low and high flows. This is as expected since  
529 the same hydrological and error models are used both for generating the “observations” and for  
530 model calibration. When real observed streamflow data is used in case  $B_2$ , there are more noticeable  
531 differences between observed and simulated streamflow. In particular, simulated peaks consistently  
532 under-estimate observed peaks. This indicates that the hydrological model and/or residual error  
533 model are miss-specified (i.e. there is evidence of “structural” model error).

534 The standardised residuals (second row in Figure 4) show large differences between the synthetic data  
535 in  $B_1$  and the real hydrological data in  $B_2$ . For  $B_1$ , standardised residuals are independent and normally  
536 distributed with zero mean and unit standard deviation. In contrast, for  $B_2$  the standardised residuals  
537 are auto-correlated, skewed (with much larger positive values than negative values), and have large  
538 extreme values ( $\sim 4$  standard deviations, c.f.  $\sim 3$  for  $B_1$ ). Regression-theory Cook’s distance depends on  
539 the magnitude of the standardised residuals (equation (5)), so these differences in standardised  
540 residuals may have a large impact on the influence metric.

541 The three leverage formulations are shown in the third row of Figure 4. Here leverage is plotted  
542 against time, rather than inputs  $\mathbf{X}$  (rainfall), so that the parabolic relationship between  $\mathbf{X}$  and linear  
543 leverage is not evident as it was in Figure 2. Linear leverage is high during rainfall events because this  
544 leverage formulation depends only on rainfall; at all other times it is zero, including immediately after  
545 these rainfall events – this is most clearly seen in Figure 4, Case  $B_1$ , column 2, row 3. In contrast,  
546 nonlinear leverage and generalised leverage remain elevated for a period of time following a rainfall  
547 event. Since generalised leverage accounts for heteroscedasticity in residual errors, it is typically

548 smaller than nonlinear leverage during high flow periods, and higher during low flow periods - this is  
549 most clearly seen in Figure 4, Case B<sub>2</sub>, column 2, row 3.

550 The magnitude of the case-deletion Cook's distance is presented in row 4 of Figure 4 as size of the  
551 grey bubbles. This influence metric is typically larger for case B<sub>2</sub> when observed streamflow is used,  
552 compared with when synthetic "observations" are used in B<sub>1</sub>. This is likely due to the impact of model  
553 mis-specification for case B<sub>2</sub>, as seen in rows 1 and 2. The accuracy of regression-theory Cook's  
554 distance compared with case-deletion Cook's distance is shown in Figure 4 (row 5). Generalised Cook's  
555 distance is the most accurate (green bubbles show the smallest absolute errors) for both cases B<sub>1</sub> and  
556 B<sub>2</sub>. For case B<sub>1</sub>, with synthetic observations, linear Cook's distance has the highest absolute errors  
557 (orange bubbles), while for case B<sub>2</sub>, real observations, nonlinear Cook's distance has the largest  
558 absolute errors.

559 Figure 5 confirms these findings when it evaluates regression-theory Cook's distance over the entire  
560 3 years of data (~1100 points). Generalised Cook's distance provides the best performance of all three  
561 regression-theory influence diagnostics, with the smallest spread about the 1:1 line and very high  
562 performance metrics (ranging from 0.93-1.00 for all metrics). Linear Cook's distance captures neither  
563 the ranking nor the values of the case-deletion Cook's distance – as reflected by the lower metrics  
564 (e.g.  $r^2$  values ranging from 0.01 to 0.23), with the sole exception of the Sp. having relatively high  
565 values (values of 0.93 and 0.90 for models B<sub>1</sub> and B<sub>2</sub>). Nonlinear Cook's distance performs a little better  
566 than linear Cook's Distance for some metrics (e.g. Sp.<sub>10</sub> improves from -0.30 to 0.95) for case B<sub>1</sub> (with  
567 synthetic observations); however, for case B<sub>2</sub> (with real observations) the performance is still relatively  
568 poor (e.g. Sp.<sub>10</sub> is 0.19 and  $r^2$  is 0.05).

569 These results indicate generalised Cook's distance is successfully able to capture the impact on  
570 leverage of the nonlinear and storage components of the hydrological model response as well as the  
571 heteroscedastic distribution of the model errors.

### 4.3. Case study set 3: Rating curve model with heteroscedastic residual errors, data uncertainty and parameter priors

The third case study set evaluates regression-theory influence diagnostics when using objective functions that account for data uncertainty and prior parameter information as part of a Bayesian inference. The magnitude of the case-deletion Cook's distance for the four rating curve cases ( $C_{1-4}$ ) are shown in Figure 6. Each panel shows observed data (with uncertainties for cases  $C_2$  and  $C_4$ ), the fitted model and the 38 case-deletion fitted models, and the relative magnitude of case-deletion Cook's distance for each data point. We provide extrapolated axes in Figure 6 to highlight the impact of influential data on the model predictions that correspond to historical evidence of the largest floods for the Ardèche River at Sauze exceeding  $6000 \text{ m}^3/\text{s}$  [Naulet et al., 2005].

In each case, the most influential data are typically extreme (both high and low) stage-discharge observed data. Accounting for discharge uncertainty in  $C_2$  (Figure 6b) slightly reduces the magnitude of the most influential data, as seen in a slight reduction of Cook's distance influence metric, and in a more practical sense in terms of reducing the variability in the case-deletion rating curves. Accounting for priors in  $C_3$  (Figure 6c) leads to a larger reduction in the influential data, while the combined effect of accounting for discharge uncertainty and priors in  $C_4$  (Figure 6d) results in an even larger reduction in the influential data, as seen by a significant reduction in case-deletion Cook's distance and a tight spread in the case-deletion rating curves. This demonstrates the value of using data uncertainty and parameter priors in reducing the impact of influential data.

Comparing the influence diagnostic results in Figure 7, the standardised residuals (second row of Figure 7) for the four rating curve models in cases  $C_{1-4}$  are quite similar, hence the leverage will largely control differences in regression-theory Cook's distance between the four cases.

The third row in Figure 7 shows the different leverage formulations for cases  $C_{1-4}$ . For linear leverage, we see the expected parabolic shape for the leverage values as a function of  $\mathbf{X}$  across the four cases

596 C<sub>1-4</sub>. As  $\mathbf{X}$  is not uniform the minima is off centre unlike the synthetic regression models case study  
597 sets (see Figure 2). For nonlinear leverage, since we have different objective functions between the  
598 cases, there are different calibrated model parameters, and hence different curves for the nonlinear  
599 leverage. Consistently the highest magnitude leverage is the highest stage-discharge value across the  
600 four cases, but the main difference in leverage occurs in the region of the knot where there is an  
601 increase in leverage as we go from C<sub>1</sub> to C<sub>2</sub>, but a decrease in leverage for C<sub>3</sub> and C<sub>4</sub>.

602 For generalised leverage there is an increase in leverage for low magnitude stage-discharge data and  
603 a decrease in leverage for high magnitude data relative to nonlinear leverage. This is because  
604 generalised leverage accounts for the heteroscedastic residual errors, which place higher weight on  
605 low vale of the stage-discharge data. There are also distinctive differences between the four cases C<sub>1-4</sub>.  
606 In C<sub>1</sub> we have higher generalised leverage than linear and nonlinear leverage with the exception of  
607 the highest stage-discharge data point where nonlinear leverage is slightly higher. Including discharge  
608 uncertainty (C<sub>2</sub>, column 2) and including prior information (C<sub>3</sub>, column 3) both result in a decrease in  
609 generalised leverage across most data points except the smallest stage measurements – with prior  
610 information especially reducing the leverage on the highest stage value. Accounting for both discharge  
611 uncertainty and priors in C<sub>4</sub> (column 3) reduces the magnitude of the generalised leverage compared  
612 to C<sub>1</sub> for all but the minimum stage measurement.

613 Figure 8 shows the performance of the three regression-theory influence diagnostics across the four  
614 rating curve models, where we see the following patterns:

- 615 1. Linear Cook's distance generally performs poorly for all data points in terms of absolute  
616 correlation ( $r^2$  range is 0.03-0.42, except for case C<sub>1</sub>) but has good performance in terms of  
617 rank correlation (Sp. range is 0.90-0.94). For the top 10 most influential points the  
618 performance is lower (Sp.<sub>10</sub> range is -0.16-0.54,  $r^2$  range is 0.01-0.33, except for C<sub>1</sub>). This  
619 indicates that the diagnostic has identified the ranking of the influential points moderately  
620 well, but does not identify the top 10 influential points.

- 621 2. Nonlinear Cook's distance has mixed performance with some mid to high range performance  
622 metrics (e.g.  $r^2$  and  $r^2_{10}$  range is 0.88-0.90 for cases  $C_1$  and  $C_2$ ) but much lower performance  
623 once the priors are incorporated (e.g.  $r^2$  and  $r^2_{10}$  range is 0.01-0.37 for cases  $C_3$  and  $C_3$ ).
- 624 3. Generalised Cook's distance has consistently high Sp. (ranging from 0.97-1.00) and performs  
625 relatively well with respect to the other metrics with lowest performance in the case of  $C_4$   
626 (Sp.<sub>10</sub> of 0.66, minimum  $r^2$  of 0.60, and minimum  $r^2_{10}$  of 0.42).

#### 627 **4.4. Performance summary of regression-theory influence diagnostics**

628 The performance metrics Sp, Sp.<sub>10</sub>,  $r^2$  and  $r^2_{10}$  for all ten cases ( $A_{1-4}$ ,  $B_{1-2}$ , and  $C_{1-4}$ ) in Sections 4.1 to 4.3  
629 are summarised in Figure 9. The results for linear Cook's distance (Figure 9, top row, columns 1 and 2)  
630 show it does a reasonable job at ranking the most influential data across all data points (very high Sp.  
631 values) However, in terms of the top-ten influential points there is a significant degradation in  
632 performance (Sp.<sub>10</sub> is lower than Sp. for all but the linear SLS model ( $A_1$ ) with some negative Sp.<sub>10</sub> for  
633 several cases meaning that the top 10 influential points are completely different to those identified  
634 by case-deletion Cook's distance. The absolute correlations (Figure 9, top row, columns 3 and 4) show  
635 that with exception of the linear SLS model, linear Cook's distance struggles to reproduce the  
636 magnitude of the case-deletion Cook's distance values.

637 Nonlinear Cook's distance (Figure 9, middle row of panels) show good performance at ranking the  
638 influential points for all data and the top 10 in synthetic cases,  $A_{1-4}$  and  $B_1$ . However for the real data  
639 case studies ( $B_2$  and  $C_{1-4}$ ) there is a sharp decrease in the performance of ranking the top 10 influential  
640 points. This is maybe because in the real case studies, the impact of the heteroscedastic residual errors  
641 comes into play, which is not accounted for by nonlinear leverage.

642 Finally we see that generalised Cook's distance (Figure 9, bottom row of panels) produces the highest  
643 performance of the regression-theory influence diagnostics across the four performance metrics. For  
644 nine of the ten case studies, all performance metrics are above 0.75. The exception being the rating

645 curve model with data uncertainty and priors ( $C_4$ ), where generalised Cook's distance, still  
646 outperforms the linear and nonlinear Cook's distance.

#### 647 **4.5. Computational efficiency of influence diagnostics**

648 An important reason for evaluating regression-theory influence diagnostics is to reduce the  
649 computational burden associated with case-deletion Cook's distance. A summary of computational  
650 demands of the different formulations is provided in Table 3, and shows that although case-deletion  
651 Cook's distance may be the most exact approach for influential point identification, it is also the most  
652 computationally intensive, requiring  $n+1$  calibration runs. In contrast, all three regression theory  
653 Cook's distance are substantially more efficient, on average requiring <1% of the computational effort  
654 of case-deletion Cook's distance.

655 Linear Cooks Distance is the fastest because regardless of the size of the calibration data set ( $n$ ) and  
656 number of model and residual error parameters ( $m$ ), it requires only one model calibration followed  
657 by the application of linear matrix algebra. Nonlinear Cook's distance has the additional computational  
658 demand of calculating the finite difference approximations for the Jacobian and Hessian matrices in  
659 the leverage formulation (equation (9)). Generalised Cook's distance has the additional computational  
660 demand of calculating the finite difference approximations for the Jacobian and Hessian matrices in  
661 the leverage formulation (equation (10)). However, surprisingly, due to the number of finite difference  
662 calculations required by each formulation, generalised leverage requires fewer model runs (~140,000  
663 in the example in Table 3) than nonlinear leverage (~270,000 runs in the example in Table 3) despite  
664 making fewer assumptions about the residual errors and therefore being broader in potential  
665 applications.

### 666 **5. Discussion**

#### 667 **5.1. Advantages and disadvantages of case-deletion and regression-theory** 668 **influence diagnostics**



669 The case-deletion and regression-theory influence diagnostics have varying assumptions and  
670 computational demands. Here we discuss the advantages and disadvantages of implementing the two  
671 classes of influence diagnostics in hydrological applications.

672 Case-deletion Cook's distance represents the most reliable measure of the influence as it provides a  
673 direct measure of the impact that a particular data point has on a model's predictions. Furthermore,  
674 hydrological models typically have nonlinear responses, including time-dependences in the  
675 predictions (and residuals) as a result of storage, and the residual errors are typically heteroscedastic  
676 and non-Gaussian. Therefore, case-deletion Cook's distance is attractive because it does not make  
677 any assumptions and can handle a wide range of modelling scenarios. However, the computational  
678 demand associated with re-calibrating the parameters for every data point in the observed record  
679 renders case-deletion influence analysis infeasible for anything but the simplest models with small  
680 datasets. For example, for a four parameter hydrological model with a decade of daily data, case-  
681 deletion required a run-time of 675 hours (~28 days) - see Table 3. A secondary concern with the  
682 implementation of case-deletion approaches is the repeated optimisation on complex response  
683 surfaces that are prone to multiple local optima [Duan et al., 1992; Kavetski et al., 2006].

684 Another drawback to applying the case-deletion Cook's distance is the loss of additional information  
685 supplied by the leverage. Cook's distance indicates which points are influential, but it does not tell us  
686 why they are influential. Analysing both the leverage and the standardised residual contribution to  
687 the magnitude of the Cook's distance therefore provides more detailed information on the nature of  
688 influential data points. Examining the standardised residuals in the case studies we see only slight  
689 variability across the four rating curve models, indicating that in some cases (such as C<sub>1-4</sub>) the leverage  
690 contribution can be the dominant factor influencing regression-theory influence diagnostics. The  
691 additional insight from examining generalised leverage is clear from a broad range of examples from  
692 the statistical literature [e.g. *Leiva et al.*, 2014; *Lemonte and Bazán*, 2015; *Osorio*, 2016; *Rocha and*  
693 *Simas*, 2011]. This is evident in the hydrological model cases B<sub>1-2</sub> where there is a clear discrepancy

694 between the magnitude of the standardised residual and the magnitude of Cook's distance, indicating  
695 the importance of the leverage in the influence of data points in the time series. In hydrological  
696 examples, points with high leverage can provide direction to the modeller in terms of where to focus  
697 additional data collection efforts. This is because these points will be highly influential in  
698 circumstances when high leverage is combined with high residual error.

699 Regression-theory influence diagnostics therefore have the following key advantages: (1) they are  
700 more efficient, due to the minimal additional computational requirements compared to a standard  
701 hydrological model calibration (99.6% fewer runs than case-deletion Cook's distance as indicated in  
702 Table 3), and (2) they provide additional diagnostic information in the form of the leverage and  
703 standardised residuals. The key limitations of regression-theory influence diagnostics are (1) they  
704 cannot evaluate case-deletion impact on predictions, parameters or objective function values (see  
705 Figure 1), and (2) they have assumptions required in the regression model structure and residual errors  
706 to formulate the leverage. In the empirical results of this study, the impact of these assumptions was  
707 illustrated with the low performance of linear and nonlinear Cook's distance on real data case studies,  
708 which had both model nonlinearity and heteroscedastic residual errors.

709 The development of generalised Cook's distance, which uses generalised leverage, to efficiently  
710 identify influential data points demonstrates considerable promise. For the ten case studies with a  
711 broad range of modelling scenarios (i.e. nonlinear model response, heteroscedastic residual error,  
712 data uncertainty and Bayesian inference) we saw generally high performance in terms of its ability to  
713 identify the same influential points as case-deletion Cook's distance at a fraction of the overall  
714 computational cost. This demonstrates that calculating generalised Cook's distance using generalised  
715 leverage provides a promising avenue to evaluate influential points in complex hydrological and  
716 environmental modelling scenarios. For future applications of influence diagnostics an attractive  
717 alternative to case-deletion and regression-theory influence diagnostics is to apply a hybrid  
718 framework for influence assessment [Wright *et al.*, 2018] that combines the strengths of the two

719 existing classes; namely 1) computational efficiency, and 2) flexibility to quantify influence using  
720 hydrologically relevant metrics.

## 721 **5.2. Application of generalised Cook's distance to a broader class of hydrological** 722 **and environmental modelling scenarios**

723 An important advantage of generalised Cook's distance is that the formulation of generalised leverage  
724 on which it is based can be applied to a very broad class of objective functions, as long they can be  
725 written in the general form in equation (2). Examples of suitable objective functions are: (1) those that  
726 account for autocorrelation in the residual error [see *Wei et al.*, 1998], which is common in  
727 hydrological modelling [see *Evin et al.*, 2014], and (2) alternative methods to account for  
728 heteroscedasticity such as logarithmic and Box-Cox transformations, also common in hydrological  
729 modelling [see *McInerney et al.*, 2017]. The additional challenges in applying generalised Cook's  
730 distance to environmental models outside of the model classes described herein could include:  
731 increased model structure complexity, increased computation time for model simulations, increased  
732 size of the parameter space, and potential challenges in numerically differentiating the objective  
733 function. A number of these challenges are in common with case-deletion approaches (e.g. the  
734 increased computational time), whereas others are unique to regression-based approaches (e.g.,  
735 numerical differentiation issues).

736 Furthermore, an extension to this work would be to examine whether removing influential data in the  
737 model calibration period can improve predictions on an independent model validation time series.  
738 This would further demonstrate the impact of influential data, given the importance of model  
739 validation in hydrology [*Biondi et al.*, 2012]

## 740 **5.3. Understanding the key drivers of influential data is key to reducing their** 741 **impact on model calibration**

742 Due to complex interactions between the chosen data, model and objective function, it can be difficult  
743 to identify influential data without undertaking an influence analysis post model calibration. Future  
744 work could endeavour to understand the key drivers of influential data by identifying situations where  
745 data are influential due to drivers independent of the choice of response model and objective function  
746 (e.g. rainfall and streamflow from an extreme weather event) and those situations where influential  
747 data are driven by the choice of response model (e.g. the response model poorly describes the  
748 response between  $y$  and  $\mathbf{X}$ ) and/or choice of objective function (e.g. the assumed residual error  
749 model poorly describes the residual error structure). Understanding these key drivers of influential  
750 data and determining whether influential data follow a particular pattern (e.g. they tend to be the  
751 largest observed model input and/or output values, or they correspond to a specific input range, etc.)  
752 will enable the modeller to determine if additional targeted data collection (e.g. collection of more  
753 high or low flows) and/or changes to the response model and/or objective function are needed to  
754 reduce the impact of influential data. The computationally efficient regression-theory influence  
755 diagnostics developed in this study will enable future investigation towards this long term goal.

## 756 **6. Conclusions**

757 Influence diagnostics identify data points that have a disproportionate impact on model parameters,  
758 performance and/or predictions, and are therefore useful tool as part of the model calibration  
759 process. Case-deletion influence diagnostics provide an exact measure of influence; however, they  
760 have a large computational demand due to the requirement for re-calibration of the model  
761 parameters for every data point in the calibration dataset. Regression-theory influence diagnostics  
762 provide an approximation of case-deletion Cook's distance by combining two regression components  
763 for each observed data point: 1) the leverage which is used to assess the potential importance of  
764 individual observations, and 2) the standardised residuals. These are more computationally efficient  
765 than case-deletion influence diagnostics, but require making assumptions about the response model  
766 and the residual error.

767 We evaluate the performance of the regression-theory influence diagnostics for three different  
768 approaches 1) linear Cook's distance, which uses linear leverage, 2) nonlinear Cook's distance, which  
769 uses nonlinear leverage, and 3) generalised Cook's distance, which uses generalised leverage. This  
770 study is the first time that generalised leverage has been combined with the standardised residual to  
771 produce generalised Cook's distance in this manner. The performance in identifying the most  
772 influential data points was evaluated against case-deletion Cook's distance on a wide range of  
773 modelling scenarios (ten case studies) that included linear/nonlinear model responses,  
774 homoscedastic/heteroscedastic residual errors, and Bayesian approaches that include data  
775 uncertainty and prior information. The performance evaluation looked at correlations (rank and  
776 absolute) with the entire dataset and the top 10 influential points identified by case-deletion Cook's  
777 distance.

778 The key outcome of this study is that generalised Cook's distance has a high performance in  
779 approximating case-deletion Cook's distance (measured by the rank and absolute correlations) for the  
780 following modelling scenarios :

- 781 1. Nonlinear regression model with heteroscedastic residual error (Sp. 0.97,  $r^2$  0.92),
- 782 2. Daily hydrological model including nonlinear model response and storage with  
783 heteroscedastic residual error (Sp. 0.93,  $r^2$  0.98),
- 784 3. Rating curve model calibrated using a Bayesian framework that includes heteroscedastic  
785 residual error, data uncertainty and prior information (Sp. 0.98,  $r^2$  0.60).

786 Importantly, generalised Cook's distance was able to achieve this high performance at identifying  
787 influential points at a fraction of the computational cost (<1%) of case-deletion Cook's distance.

788 As hydrological modelling complexity increases (i.e. more complex model structures [*Fenicia et al.*,  
789 2011], multi-catchment datasets (e.g. >200 catchments [*Coron et al.*, 2012]), and complex objective  
790 functions [*Schoups and Vrugt*, 2010]), hydrological modellers are increasingly reliant on methods to

791 detect and diagnose the impact of modelling decisions on whether a realistic representation of the  
792 catchment response has been achieved [Gupta *et al.*, 2008]. Influential data could be significant  
793 impediment towards this goal, as their presence indicates heightened sensitivity of model outputs to  
794 a small number of data points. The development of generalised Cook's distance enables influential  
795 points to be identified without the computational demand of undertaking the numerous re-  
796 calibrations required by case-deletion Cook's distance.

797 **7. References**

798 Andréassian, V., F. Bourgin, L. Oudin, T. Mathevet, C. Perrin, J. Lerat, L. Coron, and L. Berthet (2014),  
799 Seeking genericity in the selection of parameter sets: Impact on hydrological model efficiency, *Water*  
800 *Resources Research*, 50(10), 8356-8366.

801 Beven, K. (2011), *Rainfall-runoff modelling: the primer*, John Wiley & Sons.

802 Biondi, D., G. Freni, V. Iacobellis, G. Mascaro, and A. Montanari (2012), Validation of hydrological  
803 models: Conceptual basis, methodological approaches and a proposal for a code of practice, *Physics*  
804 *and Chemistry of the Earth, Parts A/B/C*, 42, 70-76.

805 Cook, R. D. (1977), Detection of Influential Observation in Linear-Regression, *Technometrics*, 19(1), 15-  
806 18.

807 Cook, R. D., and S. Weisberg (1982), *Residuals and influence in linear regression*, Chapman and Hall,  
808 New York.

809 Coron, L., V. Andréassian, C. Perrin, J. Lerat, J. Vaze, M. Bourqui, and F. Hendrickx (2012), Crash testing  
810 hydrological models in contrasted climate conditions: An experiment on 216 Australian catchments,  
811 *Water Resources Research*, 48(5).

812 Das, S. (2008), *Generalized linear models and beyond: An innovative approach from Bayesian*  
813 *perspective*, ProQuest.

814 Duan, Q. Y., S. Sorooshian, and V. Gupta (1992), Effective and efficient global optimization for  
815 conceptual rainfall-runoff models, *Water Resources Research*, 28(4), 1015-1031,  
816 doi:10.1029/1091WR02985.

817 Duan, Q. Y., S. Sorooshian, and V. K. Gupta (1994), Optimal Use of the Sce-Ua Global Optimization  
818 Method for Calibrating Watershed Models, *Journal of Hydrology*, 158(3-4), 265-284.

819 Evin, G., M. Thyer, D. Kavetski, D. McInerney, and G. Kuczera (2014), Comparison of joint versus  
820 postprocessor approaches for hydrological uncertainty estimation accounting for error  
821 autocorrelation and heteroscedasticity, *Water Resources Research*, 50(3), 2350-2375.

822 Fenicia, F., D. Kavetski, and H. H. G. Savenije (2011), Elements of a flexible approach for conceptual  
823 hydrological modeling: 1. Motivation and theoretical development, *Water Resources Research*,  
824 47(11).

825 Foglia, L., M. C. Hill, S. W. Mehl, and P. Burlando (2009), Sensitivity analysis, calibration, and testing of  
826 a distributed hydrological model using error-based weighting and one objective function, *Water*  
827 *Resources Research*, 45.

828 Foglia, L., S. W. Mehl, M. C. Hill, P. Perona, and P. Burlando (2007), Testing alternative ground water  
829 models using cross-validation and other methods, *Ground Water*, 45(5), 627-641.

830 Fox, J., and S. Weisberg (2011), *An R Companion to Applied Regression, Second Edition*, Sage  
831 Publications, Inc.

832 Gupta, H. V., T. Wagener, and Y. Liu (2008), Reconciling theory with observations: elements of a  
833 diagnostic approach to model evaluation, *Hydrological Processes*, 22(18), 3802-3813.

834 Hill, M. C., D. Kavetski, M. Clark, M. Ye, M. Arabi, D. Lu, L. Foglia, and S. Mehl (2015), Practical Use of  
835 Computationally Frugal Model Analysis Methods, *Groundwater*.

836 Hoaglin, and Welsch (1978), The Hat Matrix in Regression and ANOVA, *The American Statistician*, 32,  
837 17-22.

838 Kavetski, D., and G. Kuczera (2007), Model smoothing strategies to remove microscale discontinuities  
839 and spurious secondary optima in objective functions in hydrological calibration, *Water Resources*  
840 *Research*, 43(3).

841 Kavetski, D., G. Kuczera, and S. W. Franks (2006), Calibration of conceptual hydrological models  
842 revisited: 1. Overcoming numerical artefacts, *Journal of Hydrology*, 320(1-2), 173-186.

843 Le Coz, J., B. Renard, L. Bonnifait, F. Branger, and R. Le Boursicaud (2014), Combining hydraulic  
844 knowledge and uncertain gaugings in the estimation of hydrometric rating curves: A Bayesian  
845 approach, *Journal of Hydrology*, 509, 573-587.

846 Le Moine, N., V. Andréassian, C. Perrin, and C. Michel (2007), How can rainfall-runoff models handle  
847 intercatchment groundwater flows? Theoretical study based on 1040 French catchments, *Water*  
848 *Resources Research*, 43(6).

849 Lebecherel, L., V. Andréassian, and C. Perrin (2016), On evaluating the robustness of spatial-proximity-  
850 based regionalization methods, *Journal of Hydrology*, 539, 196-203.

851 Leiva, V., E. Rojas, M. Galea, and A. Sanhueza (2014), Diagnostics in Birnbaum-Saunders accelerated  
852 life models with an application to fatigue data, *Applied Stochastic Models in Business and Industry*,  
853 30(2), 115-131.

854 Lemonte, A. J., and J. L. Bazán (2015), New class of Johnson SB distributions and its associated  
855 regression model for rates and proportions, *Biometrical Journal*.

856 McInerney, D., M. Thyer, D. Kavetski, J. Lerat, and G. Kuczera (2017), Improving probabilistic prediction  
857 of daily streamflow by identifying Pareto optimal approaches for modeling heteroscedastic residual  
858 errors, *Water Resources Research*.

859 Nash, J. E., and J. V. Sutcliffe (1970), River flow forecasting through conceptual models part I - A  
860 discussion of principles, *Journal of Hydrology*, 10(3), 282-290.

861 Naulet, R., M. Lang, T. B. M. J. Ouarda, D. Coeur, B. Bobee, A. Recking, and D. Moussay (2005), Flood  
862 frequency analysis on the Ardeche river using French documentary sources from the last two  
863 centuries, *Journal of Hydrology*, 313(1-2), 58-78.

864 Nocedal, J., and S. J. Wright (2006), *Numerical Optimization*, Springer.

865 Osorio, F. (2016), Influence diagnostics for robust P-splines using scale mixture of normal distributions,  
866 *Annals of the Institute of Statistical Mathematics*, 68(3), 589-619.

867 Oudin, L., C. Perrin, T. Mathevet, V. Andreassian, and C. Michel (2006), Impact of biased and randomly  
868 corrupted inputs on the efficiency and the parameters of watershed models, *Journal of Hydrology*,  
869 320(1-2), 62-83.

870 Perrin, C., C. Michel, and V. Andreassian (2003), Improvement of a parsimonious model for streamflow  
871 simulation, *Journal of Hydrology*, 279(1-4), 275-289, doi:210.1016/S0022-1694(1003)00225-00227.

872 Petersen-Øverleir, A. (2004), Accounting for heteroscedasticity in rating curve estimates, *Journal of*  
873 *Hydrology*, 292(1-4), 173-181.

874 Rocha, A., and A. Simas (2011), Influence diagnostics in a general class of beta regression models, *TEST*,  
875 20(1), 95-119.

876 Schoups, G., and J. A. Vrugt (2010), A formal likelihood function for parameter and predictive inference  
877 of hydrologic models with correlated, heteroscedastic, and non-Gaussian errors, *Water Resources*  
878 *Research*, 46(10), W10531.

879 Sorooshian, S., and J. A. Dracup (1980), Stochastic parameter estimation procedures for hydrologic  
880 rainfall-runoff models: Correlated and heteroscedastic error cases, *Water Resources Research*, 16(2),  
881 430-442.

882 St. Laurent, R. T., and R. D. Cook (1992), Leverage and Superleverage in Nonlinear-Regression, *J Am*  
883 *Stat Assoc*, 87(420), 985-990.

884 St. Laurent, R. T., and R. D. Cook (1993), Leverage, local influence and curvature in nonlinear  
885 regression, *Biometrika Trust*, 80(1), 99-106

886 Thyer, M., B. Renard, D. Kavetski, G. Kuczera, S. W. Franks, and S. Srikanthan (2009), Critical evaluation  
887 of parameter consistency and predictive uncertainty in hydrological modeling: A case study using  
888 Bayesian total error analysis, *Water Resources Research*, 45(12), W00B14.

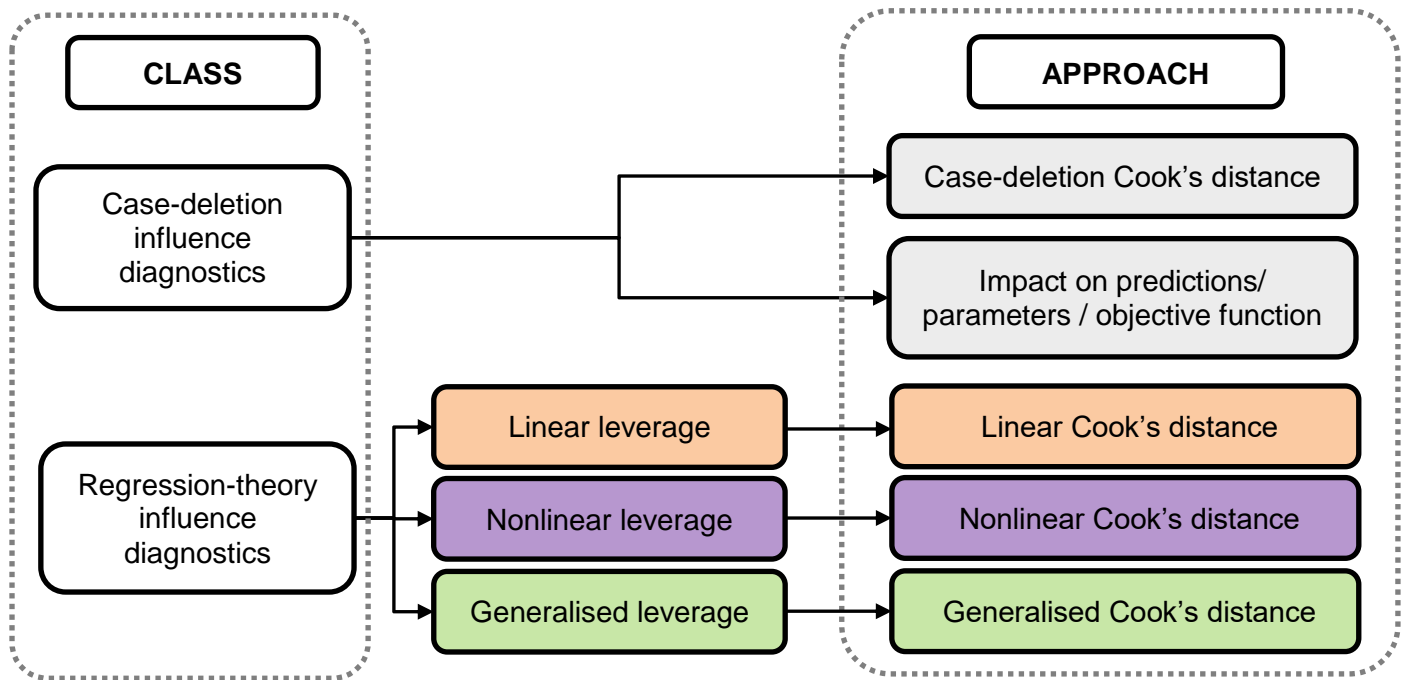
889 Wei, B. C., Y. Q. Hu, and W. K. Fung (1998), Generalized leverage and its applications, *Scandinavian*  
890 *Journal of Statistics*, 25(1), 25-37.

891 Wright, D., M. Thyer, and S. Westra (2015), Influential point detection diagnostics in the context of  
892 hydrological model calibration, *Journal of Hydrology*, 527, 1161-1172.

893 Wright, D., M. Thyer, S. Westra, and D. McInerney (2018), A hybrid framework for quantifying the  
894 influence of data in hydrological model calibration, *Journal of Hydrology*.

895

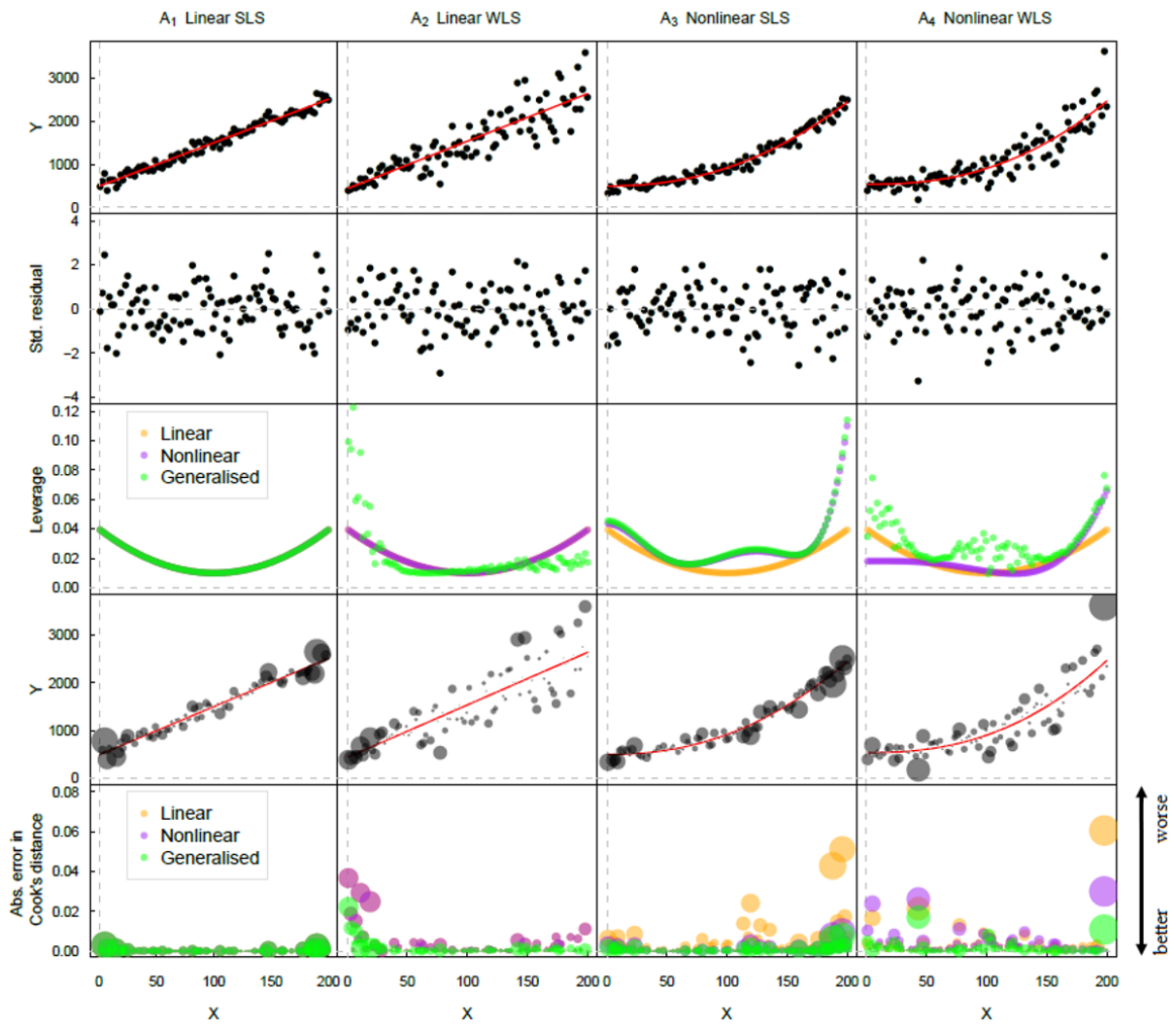




896 Figure 1 – Range of available influence diagnostics in the literature. Influence diagnostics are broken up into two classes  
 897 on the left hand side with the various approaches on the right hand side. The three regression-theory approaches are  
 898 colour coded based on the leverage formulation that they use and as they appear in the latter figures with linear Cook's  
 899 distance (orange), nonlinear Cook's distance (purple), and generalised Cook's distance (green).

900

901

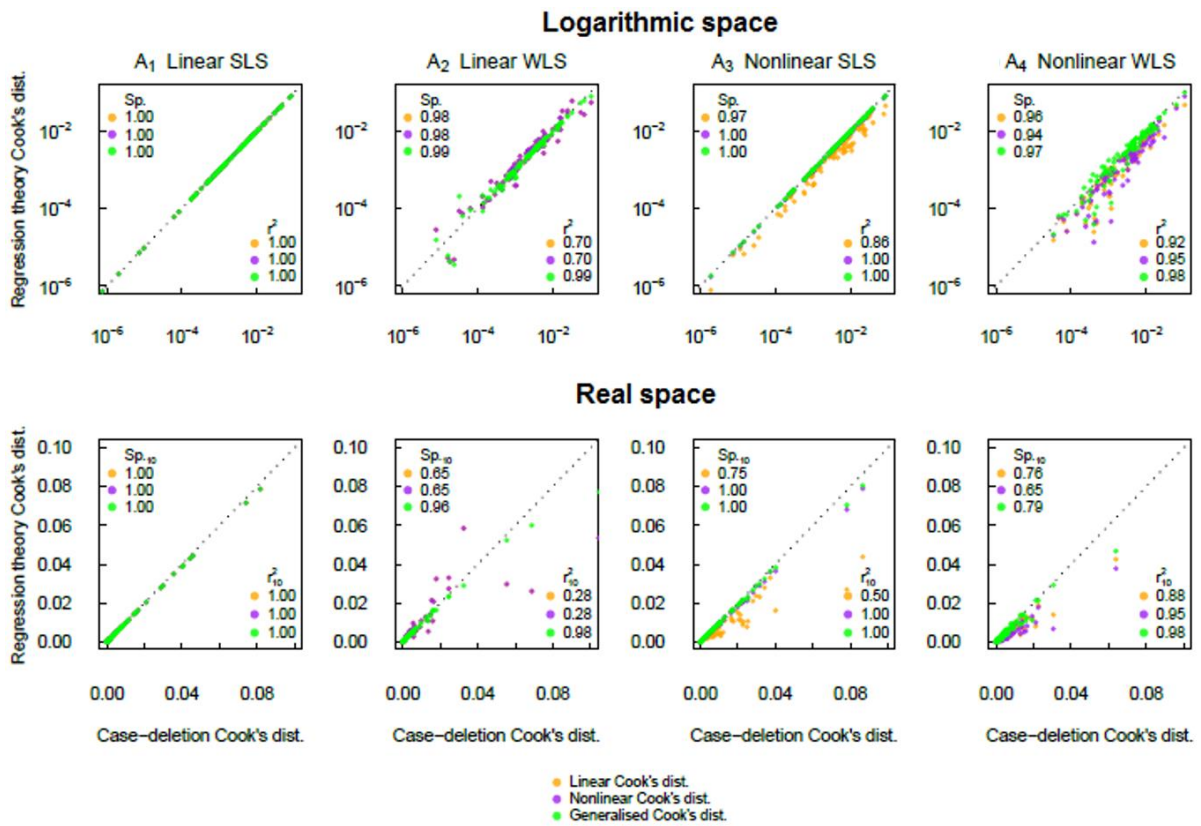


902

903 Figure 2 – Results for case study set 1: Synthetic regression models. “Observed” data (black), and model predictions (red)  
 904 in the top row, followed by standardised residuals in the second row. Leverage is shown in the third row with linear  
 905 leverage, nonlinear leverage and generalised leverage. In the case of A<sub>1</sub> the three leverage formulations are exactly equal  
 906 and so are superimposed over each other, as is the case in A<sub>2</sub> with linear and nonlinear leverage. The fourth row highlights  
 907 the distribution of the most influential data in the context of the observed data (black) and model predictions (red) where  
 908 the size of the bubbles is scaled to the value of case-deletion Cook’s distance giving a relative indication of influence. For  
 909 actual case-deletion Cook’s distance values refer to Figure 3. The final row shows the absolute error between regression-  
 910 theory Cook’s distance and case-deletion Cook’s distance where the size of the bubbles is scaled to the value of case-  
 911 deletion Cook’s distance to highlight the absolute error for the most influential data points. Note that in the final row the  
 912 relative errors are superimposed over each other.

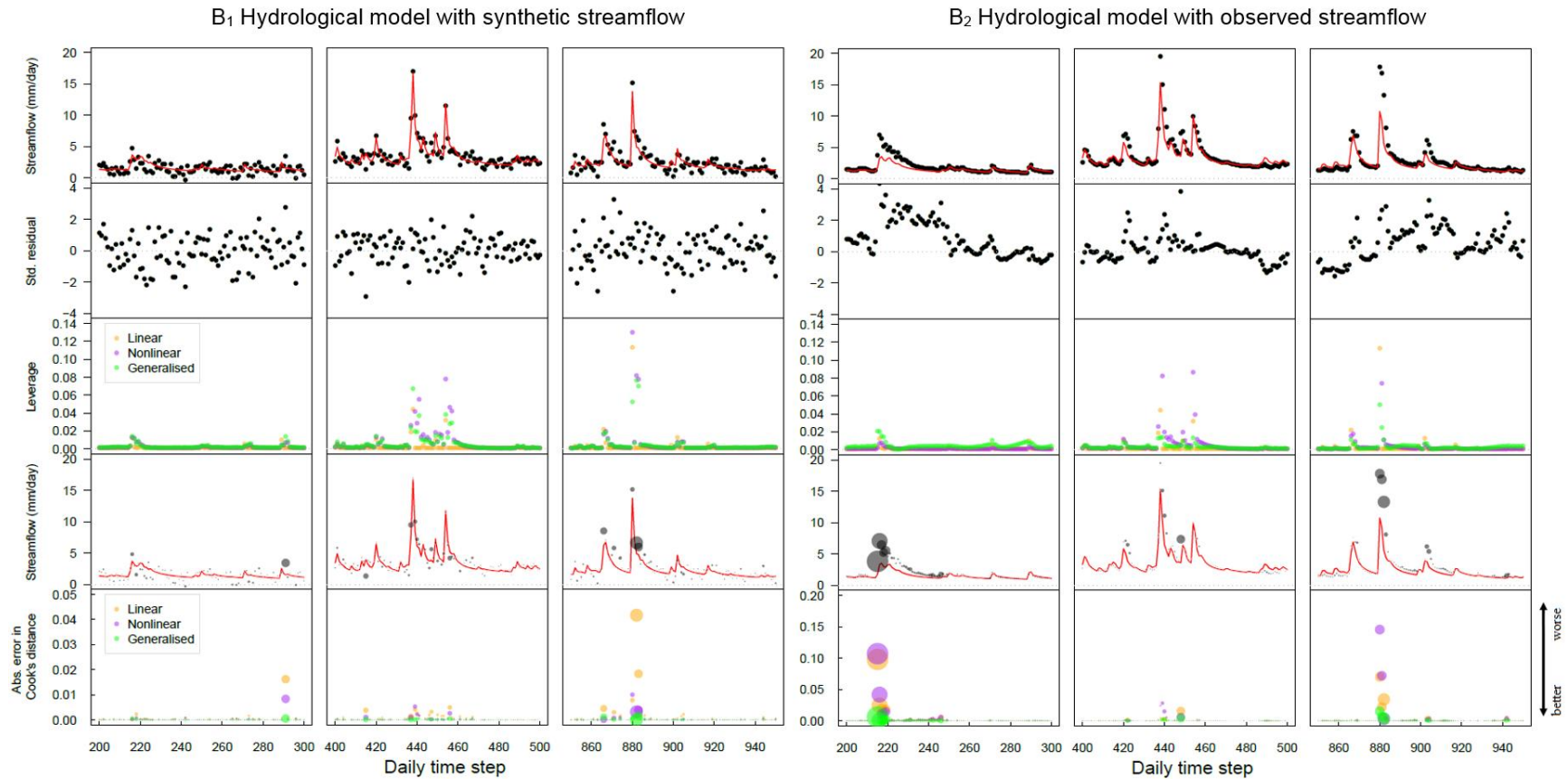
913

914



915

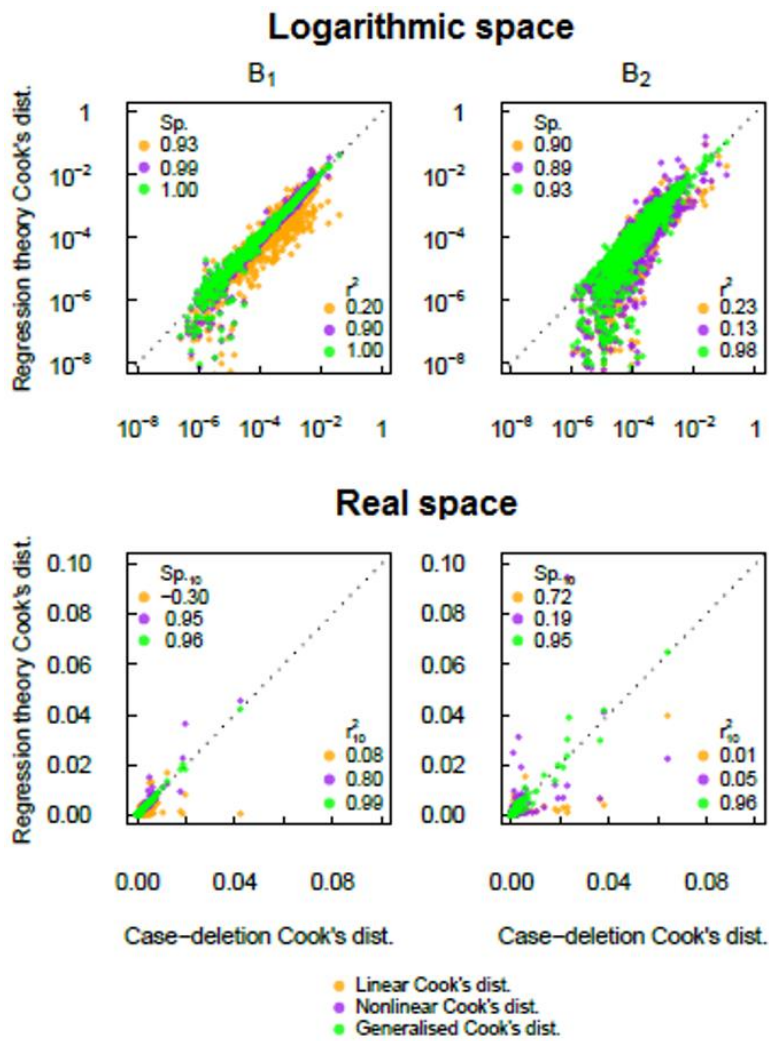
916 Figure 3 –Comparison of case-deletion Cook's distance and regression-theory influence diagnostics for case study set 1:  
 917 Synthetic regression models. In the first row we compare the performance in logarithmic space and use the Spearman  
 918 rank correlation ( $Sp$ ) and Pearson correlation ( $r^2$ ) to highlight performance across the whole dataset. In the second row  
 919 we compare the performance in real space and use the  $Sp_{.10}$  and  $r_{10}^2$  to compare the subset of the ten most influential  
 920 data points.



921

922 **Figure 4 – Results from case study set 2: Daily hydrological modelling case studies B<sub>1</sub> and B<sub>2</sub>, presented in an analogous manner to Figure 2. Observed streamflow (black), and predicted**  
 923 **streamflow (red) are shown in the top row for three different representative 100 day time periods, followed by standardised residuals in the second row. Leverage is shown in the third**  
 924 **row. The fourth row highlights the distribution of the most influential data, where the size of the bubbles is scaled to the value of case-deletion Cook's distance. The final row shows the**  
 925 **absolute error between regression-theory Cook's distance and case-deletion Cook's distance. Note that in the final row the relative errors are superimposed over each other.**

926



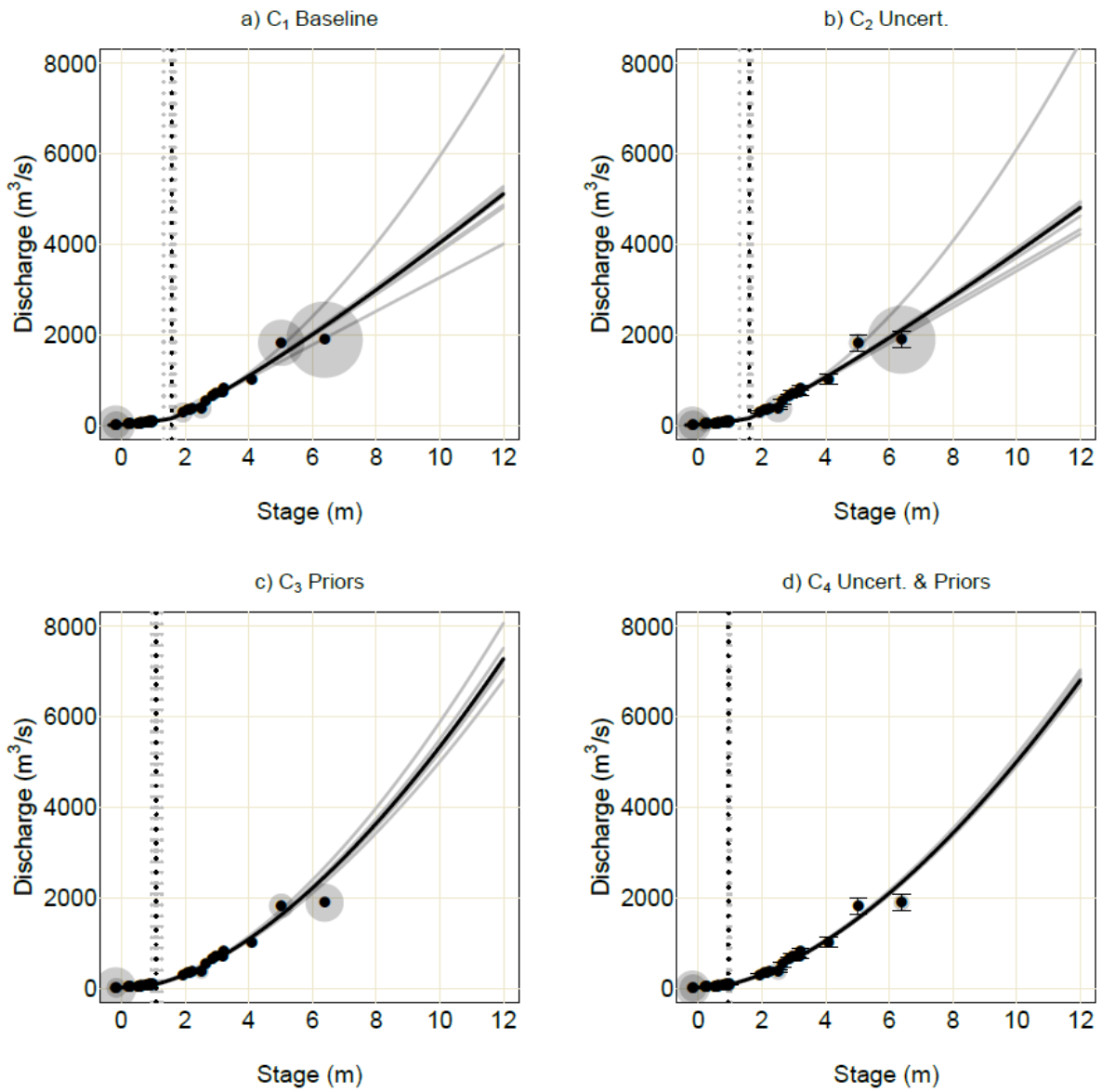
927

928 Figure 5 –Comparison of case-deletion and regression-theory influence diagnostics for case study set 2: Daily hydrological  
929 modelling cases B<sub>1</sub> and B<sub>2</sub>, presented in the same manner as Figure 3

930

931

932

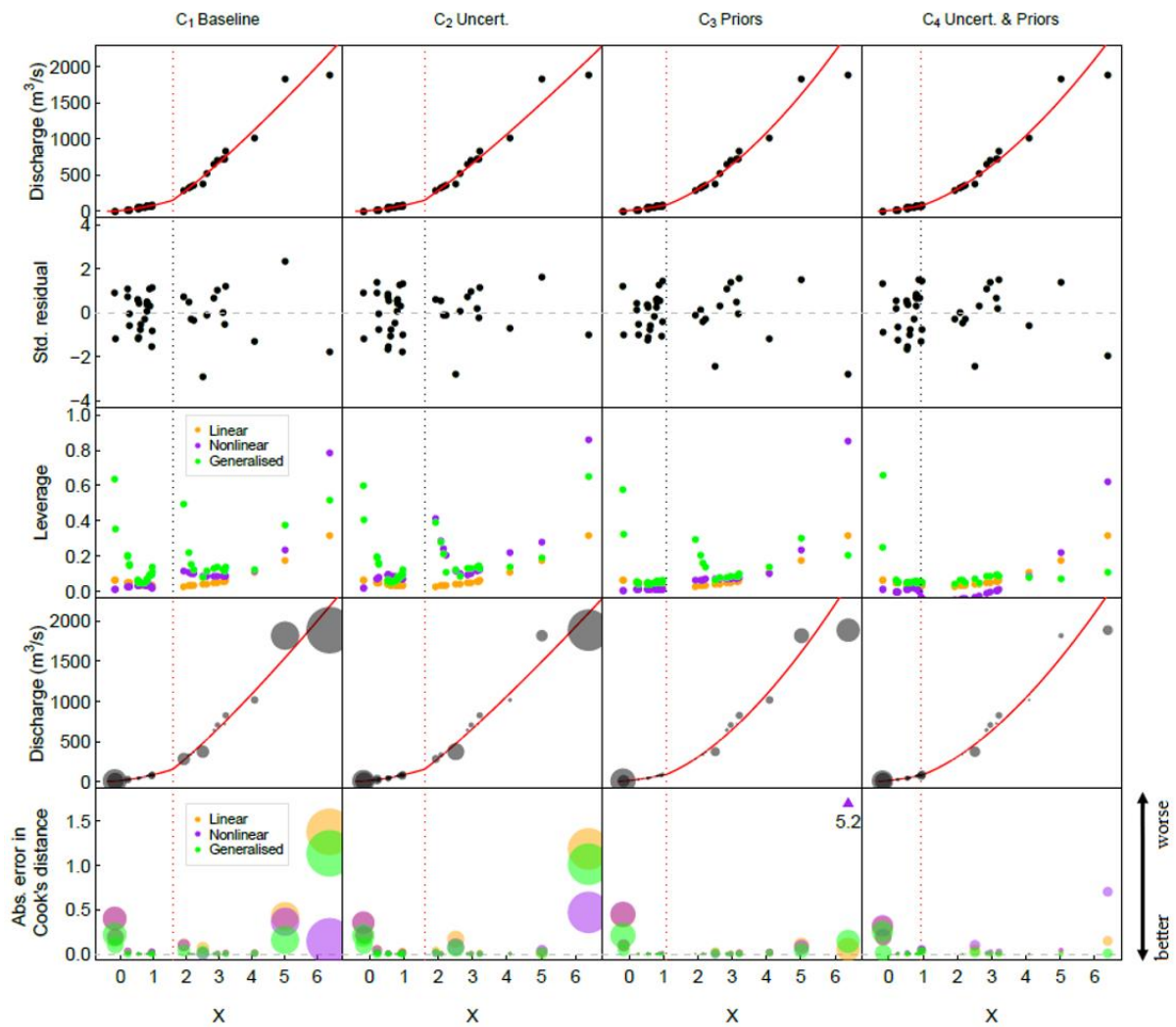


933

934 Figure 6 – Stage-discharge rating curves for the Ardèche River at Sauze. The four rating-curves presented are a) baseline  
935 rating curve without accounting for discharge uncertainty and priors, b) Rating curve with discharge uncertainty, c) Rating  
936 curve with parameter priors, d) Rating curve with both discharge uncertainty and parameter priors. Corresponding  
937 computed transition levels between section and channel controls is marked with vertical broken lines. The 38 case-  
938 deletion rating-curves and computed transition levels are shown in grey. The size of the points correspond to the relative  
939 magnitude of the case-deletion Cook's distance.

940

941



942

943 Figure 7 – Results for case study set 3: Rating curve models. The computed transition level (knot) between section and  
944 channel controls is marked with a vertical dashed line. Observed data (black), and model predictions (red) in the top row,  
945 followed by standardised residuals in the second row. Leverage is shown in the third row. The fourth row highlights the  
946 distribution of the most influential data, where the size of the bubbles is scaled to the value of case-deletion Cook's  
947 distance. The final row shows the absolute error between regression-theory Cook's distance and case-deletion Cook's  
948 distance. Note that in the final row the relative errors are superimposed over each other.

949

950

951

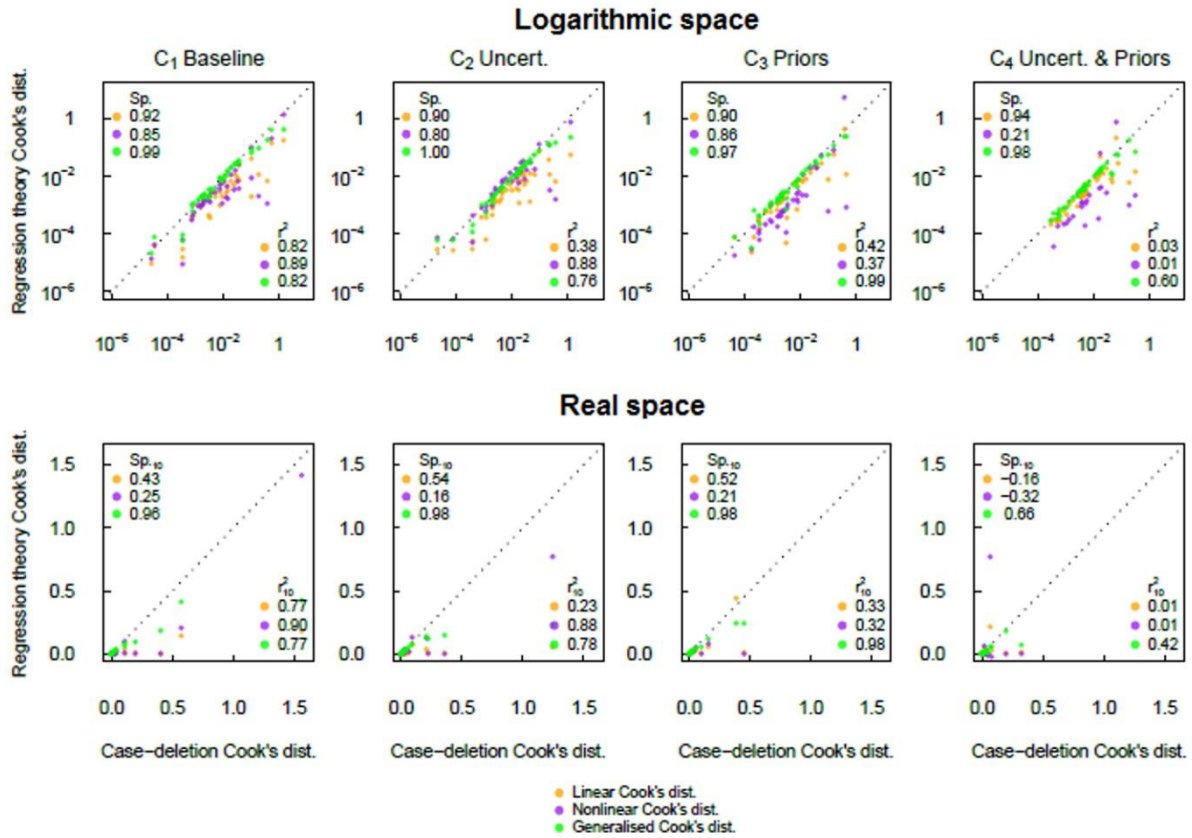
952

953



954

955



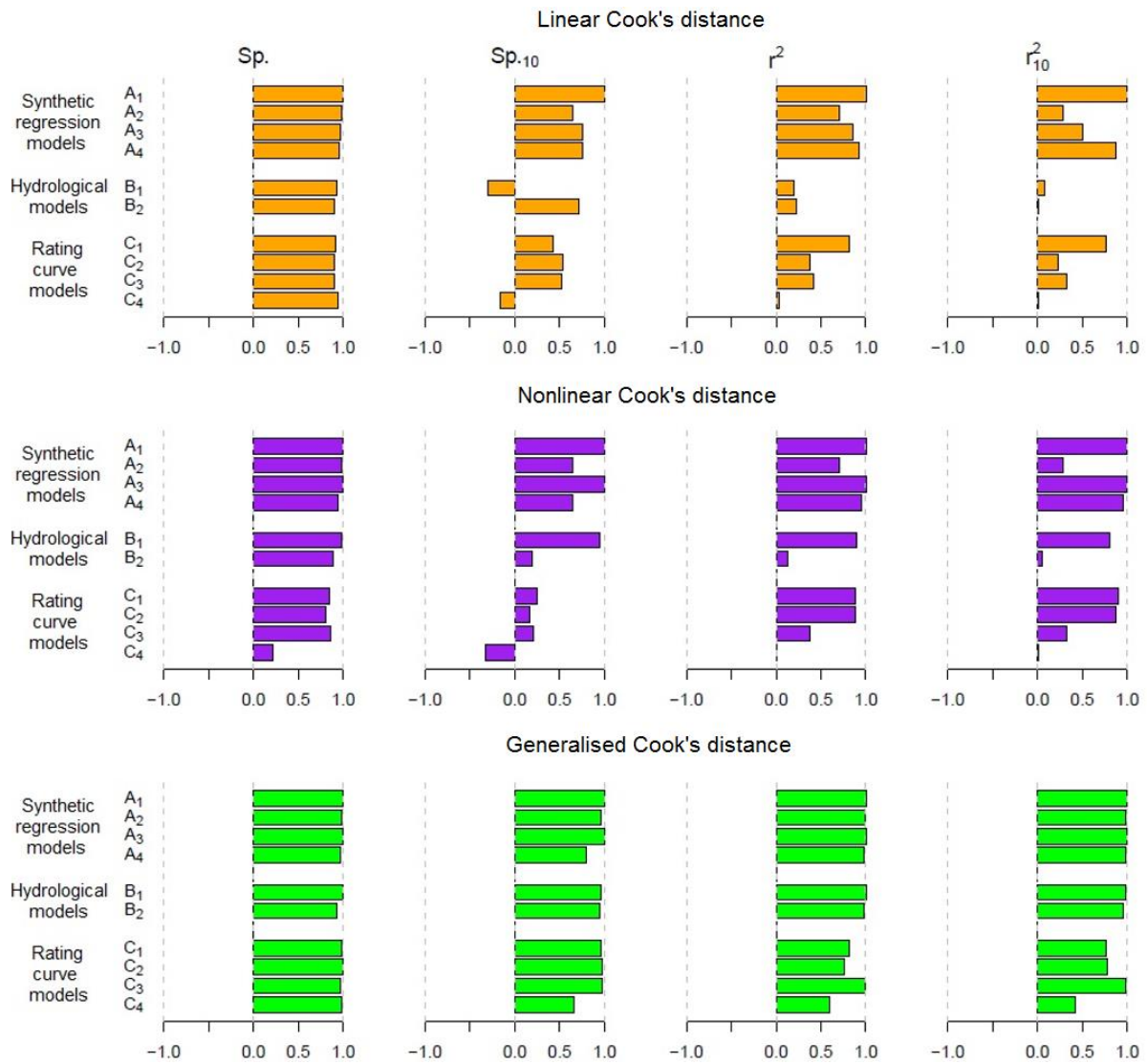
956

957 **Figure 8** Comparison of case-deletion and regression-theory influence diagnostics for case study set 3: Rating curve  
958 models, presented in the same manner as Figure 3.

959

960





961

962 **Figure 9 – Performance metrics for regression-theory influence diagnostics across the ten case studies in the three case**  
 963 **study sets. We apply the Spearman rank correlation and Pearson correlation to: (1) the whole set of data points (Sp. and**  
 964  **$r^2$ , respectively), and (2) the top 10 most influential data points identified by case-deletion Cook's distance (Sp-10 and  $r^2_{10}$ ,**  
 965 **respectfully). Linear Cook's distance is shown in the first row (orange), nonlinear Cook's distance in the second row**  
 966 **(purple) and finally generalised Cook's distance in the bottom row (green).**

967 **Table 1 – Details of the case studies.**

Case study	Response model	Residual error model	“Observed” output $Y$	Objective function
<b>Case study set 1: Synthetic regression models, Input: <math>X \sim U(0, 200)</math></b>				
<b>A<sub>1</sub></b> : Linear regression, homoscedastic residuals	$f(\mathbf{X}, \alpha_1, \alpha_2) = \alpha_1 \mathbf{X} + \alpha_2$	$\varepsilon(\sigma) \square N(0, \sigma^2), \sigma = \beta_1$	$f(\mathbf{X}, 10, 500) + \varepsilon(100)$	2.2.1
<b>A<sub>2</sub></b> : Linear regression, heteroscedastic residuals	$f(\mathbf{X}, \alpha_1, \alpha_2) = \alpha_1 \mathbf{X} + \alpha_2$	$\varepsilon(\sigma) \square N(0, \sigma^2), \sigma = \beta_1 \mathbf{y} + \beta_2$	$f(\mathbf{X}, 10, 500) + \varepsilon(0.2, 10)$	2.2.2
<b>A<sub>3</sub></b> : Nonlinear regression, homoscedastic residuals	$f(\mathbf{X}, \alpha_1, \alpha_2, \alpha_3) = \alpha_1 + \alpha_2 \mathbf{X}^{\alpha_3}$	$\varepsilon(\sigma) \square N(0, \sigma^2), \sigma = \beta_1$	$f(\mathbf{X}, 500, 0.1, 2.3) + \varepsilon(100)$	2.2.1
<b>A<sub>4</sub></b> : Nonlinear regression, heteroscedastic residuals	$f(\mathbf{X}, \alpha_1, \alpha_2, \alpha_3) = \alpha_1 + \alpha_2 \mathbf{X}^{\alpha_3}$	$\varepsilon(\sigma) \square N(0, \sigma^2), \sigma = \beta_1 \mathbf{y} + \beta_2$	$f(\mathbf{X}, 500, 0.1, 2.3) + \varepsilon(0.1, 0.5)$	2.2.2
<b>Case study set 2: Daily Hydrological models, Input: Observed rainfall measurements, All models have heteroscedastic residuals</b>				
<b>B<sub>1</sub></b> : GR4J, synthetic output	GR4J(P, PET, $\alpha$ )	$\varepsilon(\sigma) \square N(0, \sigma^2), \sigma = \beta_1 \mathbf{y} + \beta_2$	GR4J(P, PET, $\alpha = \{2200, 1.15, 87, 0.55\}$ ) + $\varepsilon(0.1, 0.5)$	2.2.2
<b>B<sub>2</sub></b> : GR4J, observed output	GR4J(P, PET, $\alpha$ )	$\varepsilon(\sigma) \square N(0, \sigma^2), \sigma = \beta_1 \mathbf{y} + \beta_2$	Observed	2.2.2
<b>Case study set 3: Rating curve models, Input: Observed stage measurements, All models have heteroscedastic residuals</b>				
<b>C<sub>1</sub></b> : Rating curve model,		$\varepsilon(\sigma) \square N(0, \sigma^2), \sigma = \beta_1 \mathbf{y} + \beta_2$		2.2.2
<b>C<sub>2</sub></b> : Rating curve model, data uncertainty		$\varepsilon(\sigma) \square N(0, \sigma^2), \sigma = \sqrt{\sigma_r^2 + \sigma_y^2}, \sigma_r = \beta_1 \mathbf{y} + \beta_2$		2.2.3
<b>C<sub>3</sub></b> : Rating curve model, parameter priors	$f(X_i, \alpha) = \begin{cases} \alpha_1 (X_i - \alpha_2)^{\alpha_3}, & X_i < \alpha_4 \\ \alpha_5 (X_i - b_2)^{\alpha_6}, & X_i \geq \alpha_4 \end{cases}$	$\varepsilon(\sigma) \square N(0, \sigma^2), \sigma = \beta_1 \mathbf{y} + \beta_2$	Observed	2.2.4
<b>C<sub>4</sub></b> : Rating curve model, data uncertainty, parameter priors		$\varepsilon(\sigma) \square N(0, \sigma^2), \sigma = \sqrt{\sigma_r^2 + \sigma_y^2}, \sigma_r = \beta_1 \mathbf{y} + \beta_2$		2.2.5

969

970

971 Table 2 – Selected prior mean (standard deviation) for the two-part rating curve model taken from Le Coz [2014]. An uninformative uniform distribution was used for the residual error  
972 model parameters. Control 1 is the rectangular sill at low flows, and Control 2 is to the rectangular channel at high flows.

	Control 1				Control 2	
$\alpha$	$a_1$	$b_1$	$c_1$	$k_1$	$a_2$	$c_2$
	50 (100)	-0.5 (2)	1.5 (0.025)	1 (1)	100(200)	1.67 (0.025)

973

974

975 Table 3 – Summary of the computational demand of case-deletion and regression-theory Cook’s distance. The example case study corresponds to the daily hydrological model (i.e.  $m_\alpha = 4$   
 976 ,  $m = 6$ ) with ~10 years of data (i.e.  $n = 3650$ ) where a fixed number of model runs is assumed per calibration ( $r = 10000$  model runs). The example runtime is calculated with a  
 977 2.90GHz processor.

Approach	Leverage	General computation demand	Model runs	Example computational demand	Example runtime (hours)	Reduction from case-deletion
<b>Case-deletion Cook’s distance</b>	-	n+1 model re-calibration	$r \times (n+1)$	36,510,000 runs	675.37	-
<b>Linear Cook’s distance</b>	Linear	Single calibration	$r$	10,000 runs	0.18	99.97%
<b>Nonlinear Cook’s distance</b>	Nonlinear	Single calibration + central difference calculations	$r + 2(n \times m_\alpha) + 4(n \times m_\alpha \times m_\alpha)$	272,800 runs	5.05	99.25%
<b>Generalised Cook’s distance</b>	Generalised	Single calibration + central difference calculations	$r + 2(n \times m) + 4(m \times m) + 4(n \times m)$	141,544 runs	2.62	99.61%

978