



## Facets of uncertainty: epistemic uncertainty, non-stationarity, likelihood, hypothesis testing, and communication

Keith Beven

To cite this article: Keith Beven (2016) Facets of uncertainty: epistemic uncertainty, non-stationarity, likelihood, hypothesis testing, and communication, Hydrological Sciences Journal, 61:9, 1652-1665, DOI: [10.1080/02626667.2015.1031761](https://doi.org/10.1080/02626667.2015.1031761)

To link to this article: <http://dx.doi.org/10.1080/02626667.2015.1031761>



© 2016 The Author(s). Published by Informa UK Limited, trading as Taylor & Francis Group.



Accepted author version posted online: 07 Apr 2015.  
Published online: 07 Jun 2016.



Submit your article to this journal [↗](#)



Article views: 1102



View related articles [↗](#)



View Crossmark data [↗](#)



Citing articles: 6 View citing articles [↗](#)

LEONARDO LECTURE



## Facets of uncertainty: epistemic uncertainty, non-stationarity, likelihood, hypothesis testing, and communication

Keith Beven<sup>a,b</sup>

<sup>a</sup>Lancaster Environment Centre, Lancaster University, Lancaster, UK; <sup>b</sup>Department of Earth Sciences, Uppsala University, Uppsala, Sweden

### ABSTRACT

This paper presents a discussion of some of the issues associated with the multiple sources of uncertainty and non-stationarity in the analysis and modelling of hydrological systems. Different forms of aleatory, epistemic, semantic, and ontological uncertainty are defined. The potential for epistemic uncertainties to induce disinformation in calibration data and arbitrary non-stationarities in model error characteristics, and surprises in predicting the future, are discussed in the context of other forms of non-stationarity. It is suggested that a condition tree is used to be explicit about the assumptions that underlie any assessment of uncertainty. This also provides an audit trail for providing evidence to decision makers.

### ARTICLE HISTORY

Received 13 February 2014  
Accepted 10 March 2015

### EDITOR

D. Koutsoyiannis

### GUEST EDITOR

S. Weijis

### KEYWORDS

Hydrological modelling;  
uncertainty estimation;  
non-stationarity; epistemic  
uncertainty; aleatory  
uncertainty; disinformation

### Introduction

I first started carrying out Monte Carlo experiments with hydrological models in 1980, while working at the University of Virginia. This was not a new approach at that time, but the computing facilities available (a CDC6600 “mainframe” computer at UVa) made it feasible for the types of hydrological model being used then. Adopting a Monte Carlo approach was a response to a personal “gut feeling” that traditional statistical approaches (at that time an analysis of uncertainty around the maximum likelihood model) were not sufficient to deal with the complex sources of uncertainty in the hydrological modelling process. Over time, we have learned much more about how to discuss facets of uncertainty in terms of aleatory, epistemic, ontological, linguistic, and other types of uncertainty (for one set of definitions see Table 1). Our perceptual model of uncertainty is now much more sophisticated but I will argue that this has not resulted in analogous progress in uncertainty quantification and, more particularly, uncertainty reduction. As one referee on this paper suggested, it can be argued that the classification of uncertainties is not really necessary: there are only epistemic uncertainties (arising from lack of knowledge) because we simply do not know enough about hydrological systems and their inputs and outputs. It is then a matter of choice as to

how to treat those uncertainties, including formal probabilistic and statistical frameworks.

What is clear is that such epistemic uncertainties will limit the inferences that can be made about hydrological systems. In particular, we are often dependent on the uncertainties associated with past observations (see, for example, Fig. 1) and have not really done a great deal about reducing hydrological data uncertainties into the past. Some observational uncertainties can certainly be treated as random variability or aleatory, but can also be subject to arbitrary uncertainties. Here, I use the word arbitrary to distinguish epistemic uncertainties that do not have simple structure or stationary statistical characteristics on the time scales used for model calibration and evaluation. This time scale qualification is important in this context since the only information we will have about the impact of different sources of uncertainties on model outputs will be contained in the sequences of model residuals within some limited period of time. It is easy to show that stochastic models based on purely aleatory variability can exhibit apparent short-period irregularity or non-stationarity (see, for example, Koutsoyiannis 2010, Montanari and Koutsoyiannis 2012). However, there is then the question of how to identify the characteristics of long-period variability from shorter periods of model residuals that might contain the type of arbitrary characteristics defined above. It has been shown that some arbitrary

**CONTACT** Keith Beven  [k.beven@lancaster.ac.uk](mailto:k.beven@lancaster.ac.uk)

© 2016 The Author(s). Published by Informa UK Limited, trading as Taylor & Francis Group.

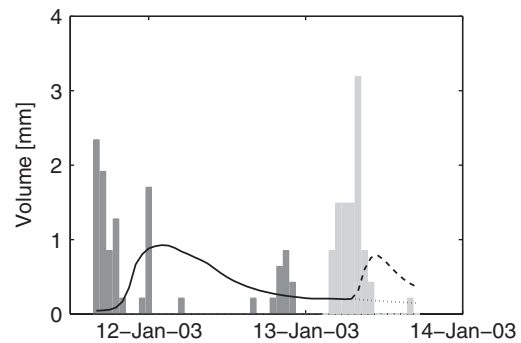
This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

**Table 1.** A classification of different types of uncertainty.

Type of uncertainty	Description
Aleatory	Uncertainty with stationary statistical characteristics. May be structured (bias, autocorrelation, long term persistence) but can be reduced to a stationary random distribution
Epistemic (system dynamics)	Uncertainty arising from a lack of knowledge about how to represent the catchment system in terms of both model structure and parameters. Note that this may include things that are included in the perceptual model of the catchment processes but are not included in the model. They may also include things that have not yet been perceived as being important but which might result in reduced model performance when surprise events occur.
Epistemic (forcing and response data)	Uncertainty arising from lack of knowledge about the forcing data or the response data with which model outputs can be evaluated. This may be because of commensurability or interpolation issues when not enough information is provided by the observational techniques to adequately describe variables required in the modelling process. May be a function of a limited gauging network, lack of knowledge about how to interpret radar data, or non-stationarity and extrapolation in rating curves.
Epistemic (disinformation)	Uncertainties in either system representation or forcing data that are <i>known</i> to be inconsistent or wrong. Real surprises. Will have the expectation of introducing disinformation into the modelling processes resulting in biased or incorrect inference (including false positives and false negatives in testing models as hypotheses).
Semantic/linguistic	Uncertainty about what statements or quantities in the relevant domain actually mean. (There are many examples in hydrology including storm runoff, baseflow, hydraulic conductivity, stationarity, etc.) This can partly result from commensurability issues that quantities with the same name have different meanings in different contexts or scales.
Ontological	Uncertainty associated with different belief systems. Relevant example here might be beliefs about whether formal probability is an appropriate framework for the representation of beliefs about the nature of model residuals. Different beliefs about the appropriate assumptions could lead to very different uncertainty estimates so that every uncertainty estimate will be conditional on the underlying beliefs and consequent assumptions.

uncertainties of this type might be *disinformative* to the model calibration process (Beven *et al.* 2011, Beven and Westerberg 2011, Kauffeldt *et al.* 2013; Fig. 1, Beven and Smith 2014), even if they might be informative in other senses (such as in identifying inconsistencies in hydrological observations, Beven and Smith 2014).

A disinformative event in this context is one for which the observational data are inconsistent with the fundamental principles (or *capacities* in the sense of Cartwright 1999) that might be applied to hydrological systems and models. Most hydrological simulation models (as opposed to forecasting models, see Beven and Young



**Figure 1.** Example of an event where the runoff coefficient based on the measured rainfalls and stream discharges is about 1.4. This clearly violates mass balance and will therefore be disinformative in calibrating a model that is constrained to maintain mass balance to represent that catchment area.

2013) impose a principle of mass balance. We expect catchment systems to also satisfy mass balance (and energy balance and momentum balance, see Reggiani *et al.* 1999). The observational data, however, might not. Figure 1 is a good example of this, with far more output as discharge from the catchment than the recorded inputs for that event. While there are some circumstances, such as a rain-on-snow event, where this could be realistic scenario, clearly no model that is constrained by mass balance would be able to reproduce such an event, suggesting that the residuals would induce bias in any model inference. It also suggests that we should take a much closer look at the data to be used in model calibration and evaluation *before* running a model (including the neglect of potential snowmelt inputs).

The implication of allowing that some model residuals might be affected by this type of arbitrary epistemic uncertainty is that commonly used probabilistic or statistical approaches to uncertainty estimation do not take enough account of the epistemic nature of uncertainty in the modelling process. It is not just a matter of finding an appropriate statistical distribution or, alternatively, some non-parametric probabilistic structure for the model residuals (e.g. Schoups and Vrugt 2010, Sikorska *et al.* 2014), especially when the sample of possible arbitrary uncertainties (or surprises) might be small. It will be suggested in what follows that we need to be more pro-active about methods for uncertainty identification and reduction. This might help to resolve some of the differences between current approaches.

### Defining types of uncertainty (and why the differences are important)

Past analysis in a variety of modelling domains in the environmental sciences has distinguished several types

of uncertainties and errors, including aleatory uncertainty, epistemic uncertainty, semantic or linguistic uncertainty, and ontological uncertainty (e.g. Beven and Binley 1992, McBratney 1992, Regan *et al.* 2002, Ascough *et al.* 2008, Beven 2009, Raadgever *et al.* 2011, Beven and Young 2013, Beven *et al.* 2014). Table 1 lists one such classification relevant to the application of hydrological models. In particular, the definition of aleatory uncertainty is constrained to the case of stationary statistical variation (noting that this might involve a structural statistical model but with stationary parameters), for which the full power of statistical theory and inference is appropriate. Epistemic uncertainties, on the other hand, have been broken down into those associated with model forcing data and observations of system response, and those associated with the representation of the system dynamics. As in Fig. 1, the observational data might sometimes be hydrologically inconsistent, and might lead to disinformation being fed into the model inference process (Beven *et al.* 2011, Beven and Smith 2014). Any of these might be sources of the rather arbitrary nature of errors in the forcing data and resulting model residual variability noted above.

Many aspects of the modelling process involve multiple sources of uncertainty, and without making very strong assumptions about the nature of these different sources it is not possible to separate the effects of the different uncertainties (Beven 2005). Attempts to separate the error associated with rainfall inputs to a catchment, for example, result in some large changes to event inputs and a strong interaction with model structural error (e.g. Vrugt *et al.* 2008, Kuczera *et al.* 2010, Renard *et al.* 2010). The very fact that there are epistemic uncertainties arising from lack of knowledge about how to represent the response, about the forcing data, and about the observed responses, reinforces this problem. If we knew what type of assumptions to make then the errors would no longer be epistemic in nature.

### Defining a method of uncertainty estimation (and why there is so much controversy about how to do so)

Uncertainty estimation has been the subject of considerable debate in the hydrological literature. There are those who consider that formal statistics is the only way to have an objective estimate of uncertainty in terms of probabilities (e.g. Mantovan and Todini 2006, Stedinger *et al.* 2008) or that the only way to deal with the unpredictable is as probabilistic variation (Montanari 2007, Montanari and Koutsoyiannis 2012). There are those who have argued that treating all

uncertainties as aleatory random variables will lead to overconfidence in model identification, so that more informal likelihood measures or limits of acceptability might be justified (e.g. within the GLUE framework of Beven 2006a, 2012, Beven and Binley 1992, 2014, Freer *et al.* 2004, Smith *et al.* 2008, Liu *et al.*, 2009; and within approximate Bayesian computation by Nott *et al.* 2012, Sadegh and Vrugt 2013, 2014). There are those who recognize the complex structure of hydrological model errors but who use transformations of different types to fit within a formal statistical framework (e.g. Montanari and Brath 2004). Some of these opinions have been explored in a number of commentaries and opinion pieces (Beven 2006a, 2006b, 2008, 2012, Hamilton 2007, Montanari 2007, Hall *et al.* 2007, Todini and Mantovan 2007, Sivakumar 2008) as well as in more technical papers.

There is, of course, no right answer—precisely because there are multiple sources of epistemic uncertainty, including model structural uncertainty, that are impossible to separate. There are also different frameworks for assessing uncertainties and different ways of formulating likelihoods. If we had knowledge of the true nature of the sources of uncertainty then they would not be epistemic and we might then be more confident about using formal statistical theory to deal with all the sources of unpredictability. Some epistemic uncertainties should be reducible by further experimentation or observation, so that there is an expectation that we might move towards more aleatory residual error in the future. In hydrology, however, this still seems a long way off, particularly with respect to the hydrological properties of the subsurface. And if, of course, there is no right answer, then this leaves plenty of scope for different philosophical and technical approaches for uncertainty estimation—or, put another way, how to define an uncertainty estimation methodology involves ontological uncertainties (Table 1). In this situation there is a lot of uncertainty about uncertainty estimation, and this is likely to be the case for the foreseeable future. This has the consequence that communication of the *meaning* of different estimates of uncertainty can be difficult. This should not, however, be an excuse for not being quite clear about the assumptions that are made in producing a particular uncertainty estimate (Faulkner *et al.* 2007, Beven and Alcock 2012, see later).

### Defining non-stationarity (in catchments and model residuals)

Many people think that the only important distinction in the modelling process is between variables that are

predictable and uncertainties that are not. Model residuals might have components of both: some identifiable predictable structure as well as some unpredictable variability. The structure indicates some aspect of the system dynamics (or boundary condition and evaluation data) that is not being captured by the model. It is often represented as a deterministic function: in the very simplest case, a stationary mean bias; in more complex cases the function might indicate some structured variability in time or space, such as a trend or seasonal component. The unpredictable component, on the other hand, is usually treated as if the variability is purely aleatory on the basis that if something is not predictable then it should be considered within a probabilistic framework (e.g. Montanari 2007) albeit that, as already noted, the nature of that variability might have some long time scale properties (Koutsoyiannis 2010, Montanari and Koutsoyiannis 2012).

This is important because it has implications for evaluating models as hypotheses in the face of epistemic errors (or long time scale aleatory errors). Hypothesis testing has traditionally been the realm of statistical inference and probability, including the recent application of Bayesian statistical theory to hydrological modelling (e.g. Clark *et al.* 2011). Purely empirically, probability and statistics can, of course, describe anything from observations to model residuals regardless of the actual sources of uncertainty as an expression of our reasonable expectations (Cox 1946). However, for any particular set of data, the resulting probabilities are conditional on the sample being considered. This is one reason why we try to abstract the empirical to a functional distributional form or the type of empirical non-parametric distributions used by Sikorska *et al.* (2014) or Beven and Smith (2014).

For simple cases where the empirical sample is random and stationary in its characteristics (after taking account of any well-defined structure) then there is a body of theory to suggest what we should expect in terms of variability in statistical characteristics as a function of sample size. There is also then a formal relationship between the statistical characteristics and a likelihood function that can be used in model evaluation. The simplest case is when the statistics of the sample have zero mean bias, constant variance, are independent and can be summarized as a Gaussian distribution. More complex likelihood functions could take account of bias, heteroscedasticity, autocorrelation, and other assumptions about the distribution. Even these more complex cases, however, are what I have called ideal cases in the past (e.g. Beven 2002, 2006a). Fundamentally, they

assume all variability in model residuals is aleatory in nature.

But real problems are not ideal in this sense; as illustrated above they are subject to arbitrary epistemic errors. It is then debatable as to whether it is appropriate to treat the errors *as if* they are aleatory. The reason is that the effective information content of any observations (or model residuals) will be reduced by epistemic uncertainties relative to the ideal case. Why is this? It is because the stationary parameter assumption of the aleatory component gives the possibility of future surprise a very low likelihood. Yet evaluating the performance of hydrological models in real applications often reveals surprises that are clearly not aleatory in this way, including occasional surprises of gross under or over predictions. This makes it difficult to define a formal statistical model of the residual structure and consequently, if the methods of estimating likelihoods in formal statistics are not valid, makes hypothesis testing of models more difficult (e.g. Beven 2010, Beven *et al.* 2012).

Consider the situation where the estimates of rainfall over a catchment might be of variable quality during a series of events in a model calibration period. The error in the estimates is not aleatory or distributional in nature because the distribution of events is not expected to be stationary (except possibly over very long periods of time, but that is not really of interest for the period of calibration data that might be available). This is the context in which we can describe the variability as rather arbitrary; i.e. we do not really know whether the rainfall uncertainties conform to any statistical distribution or if the errors in a calibration period are a good guide to the errors in the prediction period that we are actually interested in. The same could be true, of course, for aleatory errors with long-term properties (see examples in Koutsoyiannis 2010, Montanari and Koutsoyiannis 2012, Koutsoyiannis and Montanari 2015). The underlying stochastic process might then be stationary but it might be difficult to identify the properties of that process from a short-term sample with apparently non-stationary statistics. These are then both forms of epistemic uncertainty. In both cases we lack knowledge about the arbitrary nature of events or the stochastic process. We could in principle, of course, constrain that uncertainty by better observational methods, or longer data series—though that is not very useful when we only have access to calibration data collected in the past, even if we might hope to have improved data into the future.

An interesting example in this respect is the post-audit analyses of a number of groundwater modelling studies presented in Konikow and Bredehoeft (1992)

and Anderson and Woessner (1992). Model predictions of future aquifer behaviour were compared with what actually happened as the future evolved. In most studies the models failed to predict the future that actually happened. In some cases this was because, with hindsight, the original model turned out to be rather poor; in other cases it was because the future boundary conditions for the simulations had not been well predicted. In hindcasting with the correct boundary conditions the predictions were much better. Hindcasting is not all that useful, however. Where modelling is used to inform decision making (as in these groundwater cases) it is predictions of the future that are required. In these studies therefore, error characteristics were not stationary and the future turned out to hold epistemic surprises (either that the calibrated model was poor, or that the changes in boundary conditions were not those expected).

These examples involve a number of forms of non-stationarity. These are summarized in Table 2. In Class 1 we place the classical definition of non-stationarity discussed by Koutsoyiannis and Montanari (2015) in the context of stochastic process theory. They, in fact, consider that this is the *only* legitimate use of the word non-stationarity in being consistent with its technical

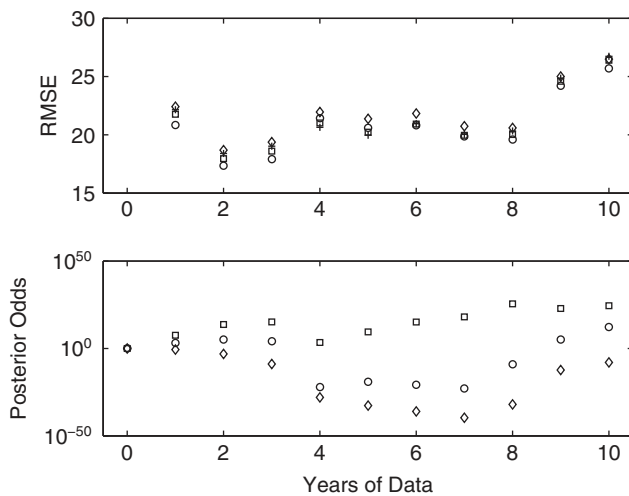
**Table 2.** Defining non-stationarity. Different classes of epistemic error that lead to non-stationarity in model residual characteristics.

Class	Source	Description
1	Non-stationarity of a stochastic process	Change over time that can be described by a deterministic function, including structure in model residuals that might compensate for consistent model or boundary condition error. All other variability will be stochastic in nature (see Koutsoyiannis and Montanari 2015).
2	Non-stationarity in catchment characteristics	Expectation that model parameters and possibly structure representing catchment characteristics will change over time or space in a way that will induce model prediction error if parameters are considered stationary
3	Non-stationarity in boundary conditions	Expectation that model boundary conditions will change over time or space in a way that will induce model prediction error if boundary conditions are poorly estimated. In some cases may include disinformative data as defined in the text.
4	Non-stationarity in model residual characteristics	Expectation that the statistical characteristics of the model residuals will vary significantly in time and space because of epistemic uncertainties about the causes of the unpredictable model error. May result from arbitrary epistemic uncertainties in boundary conditions, long-term stochastic variability, or inclusion of disinformative calibration data.

definition. In doing so, they are assuming that once any deterministic structure has been taken into account, all forms of epistemic error can be represented by a stationary stochastic model. The parameters of that model will, under the ergodic hypothesis, converge to the true values of the stochastic process as more and more observations are collected. That might, in the case of a complex stochastic process (or even some simple fractal processes) take a very large sample, but that does not negate the principle. Indeed, for a deterministic dynamical system, a stochastic representation will have stationary properties only if it is ergodic. If non-stationarity is assumed, then the system will not have ergodic properties and, Koutsoyiannis and Montanari (2015) suggest, inference will be impossible. This view means either we are back to treating all epistemic uncertainty as aleatory and stationary, once any deterministic structure has been removed, or we are simply left with unpredictability as a result of lack of knowledge.

This view has the backing of formal stochastic theory, but I think there are two issues with it. The first is the difference between what might hold in the ergodic case and the limited sample of behaviours we have in calibrating models in practical applications. The example of a stationary stochastic process giving rise to apparently non-stationary behaviour and statistics used to illustrate Koutsoyiannis and Montanari (2015) illustrates this nicely. If we have access only to a limited part of the full record, we might see periods of different statistical characteristics, or periods that include jumps. Real hydrological data might certainly be of this form, but the identification of the true stochastic process would not be possible without very long series (this is true for any fractal type behaviour). The fact that we know that the changing statistics are produced by a stationary process in such a hypothetical example, does not negate the fact that the statistics are changing and we should be wary of using an oversimplified error model (see discussion of Fig. 2 below).

Secondly, the dynamics of a nonlinear catchment model will introduce changes in the statistical properties of residuals both in the way it processes errors in the inputs and as a result of model structural error that cannot be compensated by a simple deterministic non-stationarity. From a purely hydrological point of view we expect that model residuals should have rather different characteristics on the rising limb to those around the peak to those on the falling limb in terms of bias, changing variance, and changing autocorrelation. The problem will be greater for the type of arbitrary event to event epistemic input (or model structure) error discussed above. The error in that



**Figure 2.** (Top) Root mean square errors for four model parameter sets within the same model structure (a simple single tank conceptual rainfall–runoff model, see Beven and Smith 2014). (Bottom) Likelihood ratios or posterior odds for three of the models, relative to the first (+ symbol in upper plot), evaluated using a formal likelihood and updated after the addition of further years of model residuals. The formal likelihood used allows for a mean bias, constant variance, and first-order autocorrelation and assumes a Gaussian distribution of model residuals. While similar in RMSE (and visual performance), the different models have likelihood ratios that evolve to be  $10^{40}$  different as 6 years of data are added, followed by a rapid reduction in likelihood ratio over the next 3 years.

event will also have an effect on setting up the antecedent conditions for the following event, and in some catchments, for some time into the future. The statistics of the error will be changing. Again, therefore, we should be wary of using an oversimplified error model. It is possible that again there may be a complex stochastic model that would describe all the potential changes in error statistics, but it is doubtful if it would be identifiable given the small sample of potential errors in a calibration period. It is notable that, even given a long period of calibration data, Sikorska *et al.* (2014) did not attempt to identify an underlying stochastic model of the residuals, but instead used a non-parametric probabilistic approach (in the reasonable expectation tradition of Coxian probability, Cox 1946) to represent the changing variability of the modelling uncertainties under different circumstances (see also Beven and Smith 2014). There is a difficulty with any non-parametric method, however, of how to deal with potential uncertainties in the future that are outside the range of those seen in the past.

Why is it important to make these distinctions? It is because it has an impact on what we should expect in testing a model as a hypothesis of how a catchment functions, and in particular whether it should be

considered to be fit for purpose. For example, catchments change over time (Non-stationarity Class 2) but models are often fitted with parameters that are assumed constant in time (and often space). Why is this considered acceptable practice? Perhaps, because there is an implicit expectation that this type of non-stationarity will be dominated by uncertainty in the boundary conditions used to drive a model (including the potential for Non-stationarity Class 3). There may, of course, be some clues as to whether these non-stationarities are important if there is some identifiable structure in the model residuals that could be included as a deterministic component in Non-stationarity Class 1. But we might only see the net effect of all these non-stationarities in the changing properties of the unpredictable errors (Non-stationary Class 4). But these are rarely investigated. In practical applications, statistical model inference is normally carried out *as if* all sources of error were aleatory with simple stationary properties. This assumption allows the full power of statistical inference to be applied to model calibration but would seem to be an unrealistic assumption for hydrological and other environmental models.

### Defining likelihood (and the implications for information content and hypothesis testing)

The advantage of taking a formal statistical approach to model calibration is that there is a formal link between the structure of a set of model residuals and the appropriate likelihood function. If, and only if, the assumptions about the structure of the errors are valid, then there is an additional advantage that there is a theoretical estimate of the probability of predicting a new observation. These advantages are undermined by the non-stationarities that arise from epistemic error which will generally reduce the information content (or introduce more disinformation) in the inference process than would be the case if all errors were simply aleatory with stationary parameters. So treating all sources of error as if aleatory will result in over-conditioning (and less protection against surprise in prediction). There is evidence for this in the very tight posterior parameter distributions that often arise in Bayesian calibrations of rainfall–runoff models. The likelihood surface is made very peaky such that models with very similar error variance can have tens or even hundreds of orders of magnitude difference in likelihood (Fig. 2). That really does not seem realistic to me, and did not when I first started evaluating likelihoods of multiple runs in the 1980s. The origins of the GLUE methodology lie there.

So one way ahead here might be to find more realistic likelihood functions that reflect the reduced

information content for these non-ideal cases and are robust to epistemic error. The question then is how to properly reflect the real information in a set of data when the variations are clearly not aleatory and when the summary statistics might be significantly period dependent. Again, whether the long-term properties are stationary or not is not really relevant, we want to protect against surprise in prediction (as far as is possible for an epistemic problem). In the rainfall–runoff modelling case it has been suggested that the use of summary statistics for model evaluation, such as the flow duration curve, might be more robust to error in this sense (e.g. Westerberg *et al.* 2011b, Vrugt and Sadegh 2013).

Beven *et al.* (2011) and Beven and Smith (2014) show how, for the relatively flashy South Tyne catchment in northern England (322 km<sup>2</sup>), it is possible to differentiate obviously disinformative events from informative events in model calibration within the GLUE methodology. They take an event-based approach to model evaluation that tries to reflect the relative information content expected for informative and disinformative events. They suggest that factors that will increase the relative information content of an event include: the relative accuracy of estimation of the inputs driving the model; the relative accuracy of observations with which model outputs will be compared (including commensurability issues); and the unusualness of an event (extremes, rarity of initial conditions, ...). Factors that will decrease the relative information content of an event include: repetition (multiple examples of similar conditions); inconsistency of the input and output data; the relative uncertainty of observations (e.g. highly uncertain overbank flood discharges would reduce information content of an extreme event, discharges for catchments with ill-defined rating curves might be less informative than in catchments with well defined curves); and also a preceding disinformative/less informative event over the dynamic response time scale of the catchment.

The approach depends on classifying events prior to running the model into different classes based on rainfall volume and antecedent conditions. Outlier events can be identified and examined to see if they are disinformative in terms of their runoff coefficients or other characteristics. Limits of acceptability are established for model performance in each class of informative events and a likelihood measure is based on average model performance in each class. The information content for informative events following disinformative events is weighted less highly.

Models that do not meet the limits of acceptability are rejected (given zero likelihood) in the GLUE

methodology and do not therefore contribute to the set of models to be used in prediction. This is one way of testing models as hypotheses. Epistemic error also plays a role here in that we would not want to make false negative (Type II) errors in rejecting a model that might be useful in prediction because it has been forced with poor input data. This is more serious than a false positive error in that if a poor model is not initially rejected we can hope that future evaluations would reveal its limitations. Statistical inference deals with this problem by never giving a zero likelihood, only very very small likelihoods to models that do not perform well (as seen in the orders of magnitude change in Fig. 2). This also means, however, that no model is ever rejected and hypothesis testing has to depend on some other subjective criterion, such as some informal limits on the Bayes ratios for competing models. One implication for this is that if no model is rejected there is no guarantee that the best model found is fit for purpose. This must also be assessed separately.

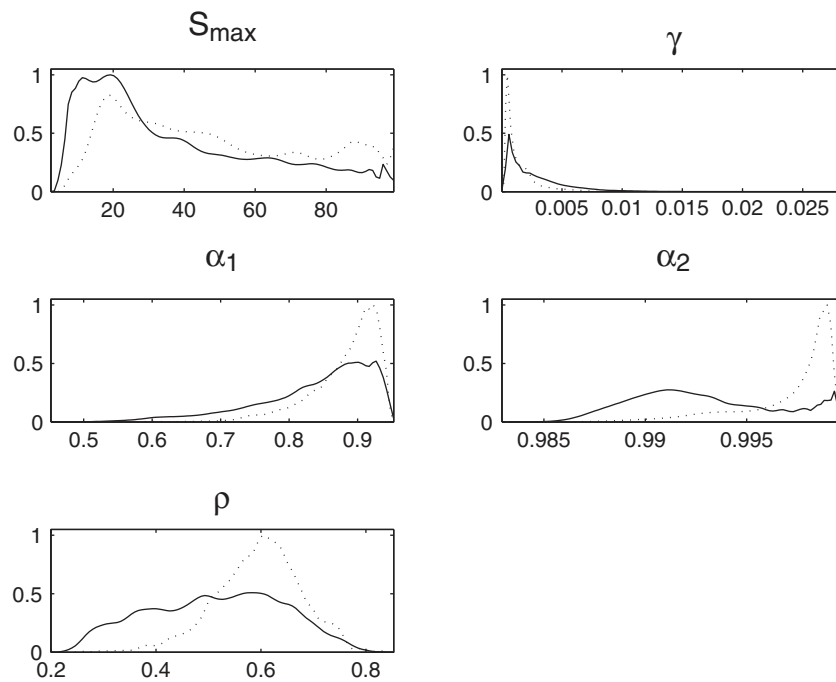
For the South Tyne catchment it turns out that using a standard dataset, as collected by the Environment Agency, there were a large number of disinformative events as distinguished by unrealistically high or low runoff coefficients. Excluding these events from the model calibration results in different posterior distributions of the model parameters (see Fig. 3). It also allows the characteristics of informative and disinformative events to be considered separately.

When it comes to prediction, however, we do not know *a priori* whether the next event will be informative or disinformative. This can only be evaluated *post hoc*, once the future has evolved (in model testing, of course, the “future” considered is some “validation” dataset). This may involve non-stationarities of error characteristics that have not been seen in the calibration period. Beven and Smith (2014) allowed for this by evaluating the error characteristics for informative and disinformative events separately and treating each new event as if it might be either informative or disinformative (Fig. 4). It was shown to help in spanning the observations for events later shown to be disinformative, but clearly cannot deal with every surprise that might occur in prediction, particularly when the system itself is non-stationary.

### Defining model rejection in hypothesis testing (and why uncertainty estimation is not the end point of a study)

In the case of the modelling study of the South Tyne catchment, some models were found that satisfied the limits of acceptability. This is not always the case; in





**Figure 3.** Posterior probability density functions for model parameters evaluated both with (solid line) and without (dotted line) calibration events classified as disinformative. Further details of this study can be found in Beven and Smith (2014).

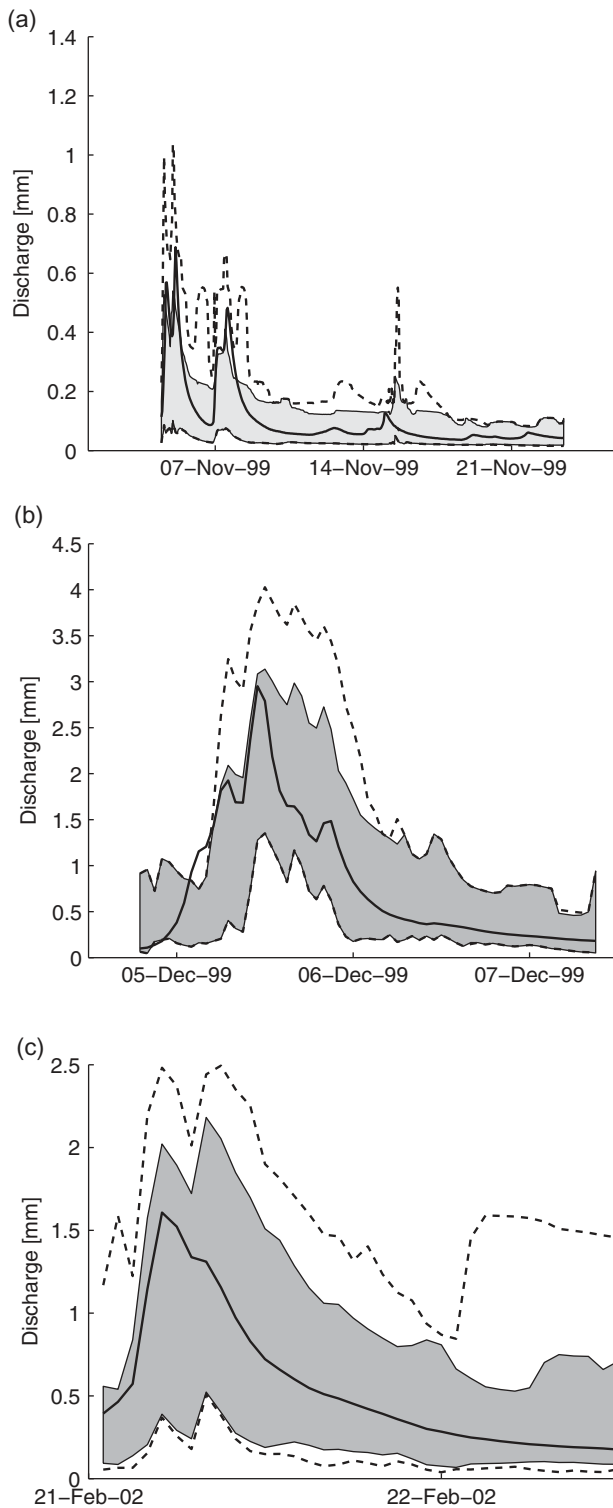
other studies no models have satisfied all the criteria of acceptability imposed (see, for example, the attempts at “blind validation” of the SHE model by Parkin *et al.* 1996, Bathurst *et al.* 2004, and the studies of Brazier *et al.* 2000, Page *et al.* 2007, Pappenberger *et al.* 2007, Choi and Beven 2007, Dean *et al.* 2009, Mitchell *et al.* 2011, within the GLUE framework using a variety of different models).

In terms of the science this is, of course, a good thing in that if all the models are rejected then improvements must be made to either the data or the model structures and parameter sets within those structures being used. That is how real progress is made. But the possibility of epistemic errors in the data used to force a model might make it difficult to make an assessment of how constrained any limits of acceptability should be. We know that all models are approximations and so such limits should be set to reflect the expectation of how well a model should be able to perform. This is a balance. We should not expect a model to predict to a greater accuracy than the assessed errors in the input and evaluation data. If it does we might suspect that it has been over-fitted to accommodate some of the particular realization of error in the calibration data.

But we also do not want to make that Type II false negative error of rejecting a model that would be useful in prediction, just because of epistemic errors and disinformation in the forcing or evaluation data.

This suggests that, if we do reject all the models tried as not fit for purpose, we should look first at the data where the model is failing and assess the potential for error in that data, especially if the failures are consistent across a large number of models. In rainfall–runoff modelling this is rarely done, but hydrological modellers are beginning to become more aware of the issues (e.g. Krueger *et al.* 2009, McMillan *et al.* 2010, 2012, Westerberg *et al.* 2011a, Kauffeldt *et al.* 2013). We also have to be careful that we have searched the model space adequately to ensure that no models have been missed. This can be difficult with high numbers of parameters, when the areas of acceptable models in the model space might be quite local. Iorgulescu *et al.* (2005) for example made 2 billion runs of a model in a 17 parameter space of which 216 were found to satisfy the (rather constrained) limits of acceptability. Blazkova and Beven (2009) made 600 000 runs of a continuous simulation flood frequency model and found that only 37 satisfied all the limits of acceptability. They also demonstrated that whether this was the case depended on the stochastic realization of the inputs used. Improved efficiency of sampling within this type of rejectionist strategy might then be valuable (e.g. the DREAM<sub>ABC</sub> code of Sadegh and Vrugt 2014).

But where all the models tried consistently fail, and we do not have any reason for suggesting that the failure is due to disinformative data, then it suggests



**Figure 4.** A sample of events taken from the model evaluation period. Each event is treated as if it is either informative (shaded 95% prediction bounds) or disinformative (dotted 95% prediction bounds). The first event is evaluated (*a posteriori*) as disinformative, the last two as informative. Further details of this study can be found in Beven and Smith (2014).

that a better model is needed. This might lead to new hypotheses about how the system is functioning, or new ways of representing some processes (see also Gupta and Nearing 2014). Model rejection is not a failure, it is an opportunity to improve either the model or data or both. Finding a better model will not provide total protection against future epistemic surprises but would, we hope, be a step in the right direction. How big a step is possible, however, will also depend on reducing uncertainty in the forcing and evaluation data.

### Communicating uncertainty to users of model predictions

There are two main reasons for incorporating uncertainty estimation into a study. One is for scientific purposes, to improve understanding of the problem and carry out hypothesis testing more rigorously. The second is because taking account of the uncertainty in model predictions might make a difference to a decision that is made in a practical application: for example, whether the planning process can take account of uncertainty in the predicted extent of flooding for the statutory design return period. For this second purpose it is necessary to communicate the *meaning* of the model predictions, and their associated uncertainties, to decision makers (e.g. Faulkner *et al.* 2007).

But, as we have seen, there can be no right answer to the estimation of uncertainty. Every estimate is conditional on the assumptions that are made and in most applications there are many assumptions that must be made (see, for example, Beven *et al.* 2014). In this case it might be useful to the communication process if the users, or particular groups of users, are introduced to the nature of those assumptions. In fact, it will generally facilitate the communication process if the users can be involved in making decisions about the relevant assumptions whenever possible. The collection of assumptions that underlie any particular application can be considered to be a form of “condition tree” (Beven and Alcock 2012, Beven *et al.* 2014). At each level of the condition tree the assumptions must be made explicit, forming an audit trail for the analysis. It has even been suggested<sup>1</sup> that every uncertainty assessment should be labelled with the names of those who produced it (and, by extension, perhaps those who agreed the assumptions on which it is based).

<sup>1</sup>For example by Jonty Rougier at Bristol University.

## Can we talk of confidence rather than uncertainty in model simulations?

Decisions about hydrological systems are made under uncertainty, and often severe uncertainty, all the time. Decision and policy makers are, however, far more interested in *evidence* than uncertainty. Evidence-based framing has become the norm in many areas of environmental policy (e.g. Boyd 2003). In the UK, the Government has considered standards for evidence (Intellectual Property Office 2011) and the Environment Agency has an Evidence Directorate and produces documents summarizing the evidence that underpins its corporate strategy. Clearly such an agency wants to have confidence in the evidence used in such policy framing. Confidence should be inversely related to error and uncertainty, but is often assessed without reference to quantifying uncertainty in either data or model results.

An example case study is the benchmarking exercise carried out to test 2D flood routing models (Environment Agency 2013). Nineteen models were tested on 12 different test cases, ranging from dam break to fluvial and urban flooding. All the test cases were hypothetical with specified roughness parameters, even if in some of the cases the geometry was based on real areas. Some had some observations available from laboratory test cases. Thus, confidence in this case represents agreement amongst models. It was shown that not all models were appropriate for all test cases, particularly those involving supercritical flow, and that some models that used simplified forms of the St. Venant equations while faster to run had more limited applicability. Differences between models depended on model implementation and numerics, so that acceptability of a model in terms of agreement with other models was essentially a subjective judgment.

There is an implicit assumption in assessing confidence in this way that in real applications to less than ideal datasets, the models that agree can be calibrated to give satisfactory simulations for mapping and planning purposes. While the report did recommend that future comparisons should also aim to assess the value of models in assessing uncertainty in the predictions, the impacts of epistemic uncertainty in defining the input, roughness parameters, and details of the geometry of the flow domain would seem to be more important than the differences between models in which we have confidence after such testing (see Beven *et al.* 2014). In real applications confidence can only be assessed by comparison with observed data, while allowing for uncertainties in inputs. Even then, there is evidence that effective values of roughness

parameters might change with the magnitude of an event, so that confidence in calibration might not carry over to more extreme events (Romanowicz and Beven 2003). Yet, for planning purposes, the Environment Agency is interested in mapping the extent of floods with annual exceedence probabilities (AEP) of 0.01 and 0.001. It is, of course, rather rare to have observations for floods within this range of AEP, more often we need to extrapolate to such levels.

It is possible to assess the uncertainty associated with such predictions and to visualize that uncertainty either as probability maps (e.g. Leedal *et al.* 2010, Neal *et al.* 2013; Beven *et al.* 2014) or as different line styles depending on the uncertainty in flood extent in different areas (Wicks *et al.* 2008). In some areas, where the flood fills the valley floor, the uncertainty in flood extent might be small, but the uncertainty in water depth, with its implications for damage calculations, might be important. In other, low slope, areas the uncertainty in extent might be significant. The advantage of doing both estimates is that confidence can be given a scale, even if, as in the Intergovernmental Panel on Climate Change (IPCC), that scale is expressed in words rather than probability. In fact, the IPCC distinguishes a scale of confidence (from “very low” to “very high”) from a scale of likelihood (from “exceptionally unlikely” to “virtually certain” based on a probability scale) (IPCC 2010). Confidence indicates how convergent the estimates of past and future change are at the current time; likelihood the degree of belief in particular future outcomes. Thus the summary of the outcomes from IPCC5 states: “Ocean warming dominates the increase in energy stored in the climate system, accounting for more than 90% of the energy accumulated between 1971 and 2010 (*high confidence*). It is *virtually certain* that the upper ocean (0–700 m) warmed from 1971 to 2010, and it *likely* warmed between the 1870s and 1971. It is *very likely* that the Arctic sea ice cover will continue to shrink and thin and that Northern Hemisphere spring snow cover will decrease during the 21st century as global mean surface temperature rises.” (IPCC 2013).

Now the IPCC will not assign any probability estimates to any of the model runs that contribute to their conclusions. They are described as projections, subject to both model limitations and conditional on scenarios of future greenhouse gas emissions. The future scenarios, and hence any probability statements, are necessarily incomplete. This has not, however, stopped the presentation of future projections in probabilistic terms in other contexts, such as those derived from an ensemble of regional model runs in the UK Climate Projections (UKCP09, see <http://ukclimateprojections>).

defra.gov.uk). The outcomes from UKCP09 are being used to assess impacts on UK hydrology (e.g. Cloke *et al.* 2010, Bell *et al.* 2012, Kay and Jones 2012) but there is sufficient epistemic uncertainty associated with both the input scenarios and the climate model implementations to be concerned about expressions of confidence or likelihood in these cases, when the probabilities may be incomplete and we should be aware of the potential for the future to surprise (Beven 2011, Wilby and Dessai 2010). Incomplete probabilities are inconsistent with a risk-based decision theoretic approach based on the exceedence probabilities of risk, although it might be possible to assess a range of exceedence curves under different assumptions about future scenarios (Rougier and Beven 2013).

We are often in this situation. Hence the need to agree assumptions and methodologies with potential users of model outcomes as discussed in the last section. Consequently, any expressions of confidence or likelihood are conditional on the assumptions, a conditionality that depends not only on what has been included, but also what might have been left out of an analysis. There will of course be epistemic uncertainties that are “unknown unknowns”. Those we do not have to worry about until, for whatever reason, they are recognized as issues and become “known unknowns”. More important are factors that are already “known unknowns”, but which are not included in the analysis because of lack of knowledge or lack of computing power or some other reason. Confidence and likelihood need to reflect the sensitivity of potential decisions to such factors since they are not easily quantified in uncertainty estimation.

### An uncertain future?

So, while quantitative uncertainty estimation is valuable in assessing the range of potential outcomes consistent with an (agreed) set of assumptions, it will generally be the case that difficult to handle epistemic uncertainties will mean that the assessment is incomplete (for good epistemic reasons). Future surprises come from that incompleteness (e.g. Beven 2013). Assessments of evidence and expressions of confidence and likelihood should reflect the potential for surprise, and robust decisions need to be insensitive to both the assessed uncertainty and the potential for surprise (erring on the side of caution, risk aversion or being precautionary). From a modeller’s perspective this has the advantage that it will reduce the possibility of a future post-audit analysis showing that the model predictions were wrong, even if why that is the case might

be obvious with hindsight (it is quite possible that this will be the case with the current generation of climate models as future improvements start to reduce the errors in predicting historical precipitation, for example).

From a decision maker’s perspective, the issues are more problematic. If, even with a detailed (and expensive) assessment of uncertainty, there remains a potential for surprise, then just how risk averse or precautionary is it necessary to be in order to make robust decisions about the future. The answer is probably that we often cannot afford to be sufficiently robust in adapting to change; it will just be too expensive. The costs and benefits of protecting against different future extremes can be assessed, even if the probability of that extreme might be difficult to estimate. In that situation, the controlling factor is likely to be the available budget (Beven 2011). That should not, of course, take away from the responsibility for ensuring that the science that underlies the evidence is as robust as possible, and communicated properly, even if those uncertainties are high and we cannot be very confident about future likelihoods in providing evidence to decision makers.

### Acknowledgements

I am grateful to Paul Smith for carrying out the model runs on which the figures are based. I am also extremely grateful to Steven Weijts, Grey Nearing, and an anonymous referee who despite disagreeing with a lot of the paper made the effort to provide very detailed comments. The paper is much improved as a consequence of their efforts, albeit that we do not agree about much that is fundamental to the issues raised.

### Disclosure statement

No potential conflict of interest was reported by the author.

### Funding

This work is a contribution to the CREDIBLE consortium funded by the UK Natural Environment Research Council (Grant NE/J017299/1). It is a written version of the Leonardo Lecture given at the Facets of Uncertainty meeting in Kos, Greece, in October 2013 with financial support from EGU.

### References

- Anderson, M. P. and Woessner, W. W., 1992. The role of the post audit in model validation. *Advances in Water Resources*, 15, 167–173. doi:10.1016/0309-1708(92)90021-S
- Ascough, J. C., *et al.*, 2008. Future research challenges for incorporation of uncertainty in environmental and

- ecological decision-making. *Ecological Modelling*, 219 (3–4), 383–399. doi:10.1016/j.ecolmodel.2008.07.015
- Bathurst, J. C., et al., 2004. Validation of catchment models for predicting land-use and climate change impacts. 3. Blind validation for internal and outlet responses. *Journal of Hydrology*, 287 (1–4), 74–94. doi:10.1016/j.jhydrol.2003.09.021
- Bell, V.A., et al., 2012. How might climate change affect river flows across the Thames Basin? An area-wide analysis using the UKCP09 Regional Climate Model ensemble. *Journal of Hydrology*, 442, 89–104. doi:10.1016/j.jhydrol.2012.04.001
- Beven, K.J., 2002. Towards a coherent philosophy for modelling the environment. *Proceedings of the Royal Society, London A*, 458, 2465–2484. doi:10.1098/rspa.2002.0986
- Beven, K.J., 2005. On the concept of model structural error. *Water Science and Technology*, 52 (6), 165–175.
- Beven, K.J., 2006a. A manifesto for the equifinality thesis. *Journal of Hydrology*, 320, 18–36. doi:10.1016/j.jhydrol.2005.07.007
- Beven, K.J., 2006b. On undermining the science? *Hydrological Processes*, 20, 3141–3146. doi:10.1002/hyp.6396
- Beven, K.J., 2008. On doing better hydrological science. *Hydrological Processes*, 22, 3549–3553. doi:10.1002/hyp.7108
- Beven, K.J., 2009. *Environmental modelling: an uncertain future?* London: Routledge.
- Beven, K.J., 2010. Preferential flows and travel time distributions: defining adequate hypothesis tests for hydrological process models. *Hydrological Processes*, 24, 1537–1547. doi:10.1002/hyp.7718
- Beven, K.J., 2011. I believe in climate change but how precautionary do we need to be in planning for the future? *Hydrological Processes*, 25, 1517–1520. doi:10.1002/hyp.7939
- Beven, K. J., 2012. Causal models as multiple working hypotheses about environmental processes. *Comptes Rendus Geoscience, Académie de Sciences*, 344, 77–88. Paris. doi:10.1016/j.crte.2012.01.005
- Beven, K.J., 2013. So how much of your error is epistemic? Lessons from Japan and Italy. *Hydrological Processes*, 27 (11), 1677–1680. doi:10.1002/hyp.9648
- Beven, K.J. and Alcock, R., 2012. Modelling everything everywhere: a new approach to decision making for water management under uncertainty. *Freshwater Biology*, 56, 10.1111/j.1365-2427.2011.02592.x
- Beven, K.J. and Binley, A.M., 1992. The future of distributed models: model calibration and uncertainty prediction. *Hydrological Processes*, 6, 279–298. doi:10.1002/hyp.3360060305
- Beven, K.J. and Binley, A. M., 2014. GLUE: 20 years on. *Hydrological Processes*, 28 (24), 5897–5918. doi:10.1002/hyp.10082
- Beven, K.J., Leedal, D. T., and McCarthy, S., 2014. Framework for assessing uncertainty in fluvial flood risk mapping. CIRIA report C721 2014. Available from [http://www.ciria.org/Resources/Free\\_publications/fluvial\\_flood\\_risk\\_mapping.aspx](http://www.ciria.org/Resources/Free_publications/fluvial_flood_risk_mapping.aspx)
- Beven, K.J., et al., 2012. Comment on “Pursuing the method of multiple working hypotheses for hydrological modeling” by P. Clark et al. *Water Resources Research*, 48, W11801. doi:10.1029/2012WR012282
- Beven, K.J. and Smith, P. J., 2014. Concepts of information content and likelihood in parameter calibration for hydrological simulation models. *ASCE Journal of Hydrologic Engineering*. doi:10.1061/(ASCE)HE.1943-5584.0000991
- Beven, K.J., Smith, P. J., and Wood, A., 2011. On the colour and spin of epistemic error (and what we might do about it). *Hydrology and Earth System Sciences*, 15, 3123–3133. doi:10.5194/hess-15-3123-2011
- Beven, K.J. and Westerberg, I., 2011. On red herrings and real herrings: disinformation and information in hydrological inference. *Hydrological Processes*, 25 (10), 1676–1680. doi:10.1002/hyp.7963
- Beven, K.J. and Young, P., 2013. A guide to good practice in modeling semantics for authors and referees. *Water Resources Research*, 49, 5092–5098. doi:10.1002/wrcr.20393
- Blazkova, S. and Beven, K.J., 2009. A limits of acceptability approach to model evaluation and uncertainty estimation in flood frequency estimation by continuous simulation: Skalka catchment, Czech Republic. *Water Resources Research*, 45, W00B16. doi:10.1029/2007WR006726
- Boyd, I., 2003. Making science count in government. *Elife*, 2, e01061. doi:10.7554/eLife.01061
- Brazier, R. E., et al., 2000. Equifinality and uncertainty in physically-based soil erosion models: application of the GLUE methodology to WEPP, the Water Erosion Prediction Project – for sites in the UK and USA. *Earth Surface Processes and Landforms*, 25, 825–845. doi:10.1002/1096-9837(200008)25:8<825::AID-ESP101>3.0.CO;2-3
- Cartwright, N., 1999. *The dappled world: a study of the boundaries of science*. Cambridge, UK: Cambridge University Press.
- Choi, H. T. and Beven, K. J., 2007. Multi-period and multi-criteria model conditioning to reduce prediction uncertainty in an application of TOPMODEL within the GLUE framework. *Journal of Hydrology*, 332 (3–4), 316–336. doi:10.1016/j.jhydrol.2006.07.012
- Clark, M.P., Kavetski, D., and Fenicia, F., 2011. Pursuing the method of multiple working hypotheses for hydrological modeling. *Water Resources Research*, 47, W09301. doi:10.1029/2010WR009827
- Cloke, H. L., et al., 2010. Climate impacts on river flow: projections for the Medway catchment, UK, with UKCP09 and CATCHMOD. *Hydrological Processes*, 24, 3476–3489. doi:10.1002/hyp.7769
- Cox, R. T., 1946. Probability, frequency and reasonable expectation. *American Journal of Physics*, 14, 1–13. doi:10.1119/1.1990764
- Dean, S., et al., 2009. Uncertainty assessment of a process-based integrated catchment model of phosphorus (INCA-P). *Stochastic Environmental Research and Risk Assessment*, 23, 991–1010. doi:10.1007/s00477-008-0273-z
- Environment Agency, 2013. Benchmarking the latest generation of 2D hydraulic modelling packages. Report SC120002, Environmental Agency: Bristol.
- Faulkner, H., et al., 2007. Developing a translational discourse to communicate uncertainty in flood risk between science and the practitioner. *AMBIO: A Journal of the Human Environment*, 36 (8), 692–704. doi:10.1579/0044-7447(2007)36[692:DATDTC]2.0.CO;2

- Freer, J.E., *et al.*, 2004. Constraining dynamic TOPMODEL responses for imprecise water table information using fuzzy rule based performance measures. *Journal of Hydrology*, 291, 254–277. doi:10.1016/j.jhydrol.2003.12.037
- Gupta, H. V. and Nearing, G. S., 2014. Debates—the future of hydrological sciences: A (common) path forward? Using models and data to learn: A systems theoretic perspective on the future of hydrological science. *Water Resources Research*, 50, 5351–5359. doi:10.1002/2013WR015096
- Hall, J, O'Connell, E, and Ewen, J., 2007. On not undermining the science: coherence, validation and expertise. Discussion of invited commentary by Keith Beven Hydrological Processes, 20, 3141–3146 (2006). *Hydrological Processes*, 21 (7), 985–988. doi:10.1002/hyp.6639
- Hamilton, S., 2007. Just say NO to equifinality. *Hydrological Processes*, 21 (14), 1979–1980. doi:10.1002/hyp.6800
- Intellectual Property Office, 2011. Good evidence for policy, UK Government. Available from <http://www.ipo.gov.uk/consult-2011-copyright-evidence.pdf>
- Iorgulescu, I, Beven, K J, and Musy, A, 2005. Data-based modelling of runoff and chemical tracer concentrations in the Haute-Mentue research catchment (Switzerland). *Hydrological Processes*, 19, 2557–2573. doi:10.1002/hyp.5731
- IPCC, 2010. Guidance note for lead authors of the 5th IPCC Assessment on Consistent Treatment of Uncertainties. Available from: <http://www.ipcc.ch/pdf/supporting-material/uncertainty-guidance-note.pdf>
- IPCC, 2013. Headline statements from the summary for policymakers. Available from: [http://www.ipcc.ch/news\\_and\\_events/docs/ar5/ar5\\_wg1\\_headlines.pdf](http://www.ipcc.ch/news_and_events/docs/ar5/ar5_wg1_headlines.pdf)
- Jain, A. and Dubes, R., 1998. *Algorithms for clustering data*. Upper Saddle River, NJ: Prentice Hall.
- Kauffeldt, A., *et al.*, 2013. Disinformative data in large-scale hydrological modelling. *Hydrology and Earth System Sciences*, 17, 2845–2857. doi:10.5194/hess-17-2845-2013
- Kay, A. L. and Jones, R. G., 2012. Comparison of the use of alternative UKCP09 products for modelling the impacts of climate change on flood frequency. *Climatic Change*, 114 (2), 211–230. doi:10.1007/s10584-011-0395-z
- Konikow, L. F. and Bredehoeft, J. D., 1992. Groundwater models cannot be validated? *Advances in Water Resources*, 15, 75–83. doi:10.1016/0309-1708(92)90033-X
- Koutsoyiannis, D., 2010. HESS opinions: “A random walk on water”. *Hydrology and Earth System Sciences*, 14 (3), 585–601. doi:10.5194/hess-14-585-2010
- Koutsoyiannis, D. and Montanari, A., 2015. Negligent killing of scientific concepts: the stationarity case. *Hydrological Sciences Journal*, 60, 7–8. doi:10.1080/02626667.2014.959959
- Krueger, T., *et al.*, 2009. Uncertainties in data and models to describe event dynamics of agricultural sediment and phosphorus transfer. *Journal of Environmental Quality*, 38 (3), 1137–1148. doi:10.2134/jeq2008.0179
- Kuczera, G., *et al.*, 2010. There are no hydrological monsters, just models and observations with large uncertainties! *Hydrological Sciences Journal*, 55 (6), 980–991. doi:10.1080/02626667.2010.504677
- Leedal, D. T., *et al.*, 2010. Visualization approaches for communicating real-time flood forecasting level and inundation information. *Journal of Flood Risk Management*, 3, 140–150. doi:10.1111/j.1753-318X.2010.01063.x
- Liu, Y-L., *et al.*, 2009. Towards a limits of acceptability approach to the calibration of hydrological models: extending observation error. *Journal of Hydrology*, 367 (1–2), 93–103. doi:10.1016/j.jhydrol.2009.01.016
- Mantovan, P. and Todini, E., 2006. Hydrological forecasting uncertainty assessment: incoherence of the GLUE methodology. *Journal of Hydrology*, 330, 368–381. doi:10.1016/j.jhydrol.2006.04.046
- McBratney, A. B., 1992. On variation, uncertainty and informatics in environmental soil management. *Australian Journal of Soil Research*, 30 (6), 913–935. doi:10.1071/SR9920913
- McMillan, H., *et al.*, 2010. Impacts of uncertain river flow data on rainfall-runoff model calibration and discharge predictions. *Hydrological Processes*, 24 (10), 1270–1284.
- McMillan, H., Krueger, T., and Freer, J., 2012. Benchmarking observational uncertainties for hydrology: rainfall, river discharge and water quality. *Hydrological Processes*, 26, 4078–4111. doi:10.1002/hyp.9384
- Mitchell, S., *et al.*, 2011. Processes influencing model-data mismatch in drought-stressed, fire-disturbed, eddy flux sites. *JGR-Biosciences*, 116, 10.1029/2009JG001146
- Montanari, A., 2007. What do we mean by ‘uncertainty’? The need for a consistent wording about uncertainty assessment in hydrology. *Hydrological Processes*, 21, 841–845. doi:10.1002/hyp.6623
- Montanari, A. and Brath, A., 2004. A stochastic approach for assessing the uncertainty of rainfall-runoff simulations. *Water Resources Research*, 40, W01106. doi:10.1029/2003WR002540
- Montanari, A. and Koutsoyiannis, D., 2012. A blueprint for process-based modeling of uncertain hydrological systems. *Water Resources Research*, 48, W09555. doi:10.1029/2011WR011412
- Neal, J., *et al.*, 2013. Probabilistic flood risk mapping including spatial dependence. *Hydrological Processes*, 27, 1349–1363. doi:10.1002/hyp.9572
- Nott, D. J., Marshall, L., and Brown, J., 2012. Generalized likelihood uncertainty estimation (GLUE) and approximate Bayesian computation: what's the connection? *Water Resources Research*, 48, 10.1029/2011WR011128
- Page, T., *et al.*, 2007. Modelling the chloride signal at Plynlimon, Wales, using a modified dynamic TOPMODEL incorporating conservative chemical mixing (with uncertainty). *Hydrological Processes*, 21, 292–307. doi:10.1002/hyp.6186
- Pappenberger, F., *et al.*, 2007. Grasping the unavoidable subjectivity in calibration of flood inundation models: a vulnerability weighted approach. *Journal of Hydrology*, 333, 275–287. doi:10.1016/j.jhydrol.2006.08.017
- Parkin, G., *et al.*, 1996. Validation of catchment models for predicting land-use and climate change impacts. 2. Case study for a Mediterranean catchment. *Journal of Hydrology*, 175 (1–4), 595–613. doi:10.1016/S0022-1694(96)80027-8
- Raadgever, G. T., *et al.*, 2011. Uncertainty management strategies: lessons from the regional implementation of the Water Framework Directive in the Netherlands. *Environmental Science & Policy*, 14 (1), 64–75. doi:10.1016/j.envsci.2010.11.001
- Regan, H. M., Colyvan, M., and Burgman, M. A., 2002. A taxonomy and treatment of uncertainty for ecology and conservation biology. *Ecological Applications*, 12 (2),

- 618–628. doi:10.1890/1051-0761(2002)012[0618:ATATOU]2.0.CO;2
- Reggiani, P., *et al.*, 1999. A unifying framework for watershed thermodynamics: constitutive relationships. *Advances in Water Resources*, 23, 15–39. doi:10.1016/S0309-1708(99)00005-6
- Renard, B., *et al.*, 2010. Understanding predictive uncertainty in hydrologic modeling: the challenge of identifying input and structural errors. *Water Resources Research*, 46, W05521, 22pp. doi:10.1029/2009WR008328
- Romanowicz, R. and Beven, K. J., 2003. Bayesian estimation of flood inundation probabilities as conditioned on event inundation maps. *Water Resources Research*, 39 (3), W01073. doi:10.1029/2001WR001056
- Rougier, J. and Beven, K. J., 2013. Model limitations: the sources and implications of epistemic uncertainty. In: J. Rougier, S. Sparks, and L. Hill, eds. *Risk and uncertainty assessment for natural hazards*. Cambridge, UK: Cambridge University Press, 40–63.
- Sadegh, M. and Vrugt, J. A., 2013. Bridging the gap between GLUE and formal statistical approaches: approximate Bayesian computation. *Hydrology and Earth System Sciences*, 17, 4831–4850. doi:10.5194/hess-17-4831-2013
- Sadegh, M. and Vrugt, J. A., 2014. Approximate Bayesian Computation using Markov Chain Monte Carlo simulation: DREAM(ABC). *Water Resources Research*, 50, 6767–6787. doi:10.1002/2014WR015386
- Schoups, G. and Vrugt, J. A., 2010. A formal likelihood function for parameter and predictive inference of hydrologic models with correlated, heteroscedastic and non-Gaussian errors. *Water Resources Research*, 46, W10531. doi:10.1029/2009WR008933
- Sikorska, A. E., Montanari, A., and Koutsoyiannis, D., 2014. Estimating the uncertainty of hydrological predictions through data-driven resampling techniques. *ASCE Journal of Hydrology Engineering*, doi:10.1061/(ASCE)HE.1943-5584.0000926
- Sivakumar, B., 2008. Undermining the science or undermining Nature? *Hydrological Processes*, 22, 893–897. doi:10.1002/hyp.7004
- Smith, P. J., Beven, K. J., and Tawn, J. A., 2008. Informal likelihood measures in model assessment: theoretic development and investigation. *Advances in Water Resources*, 31, 1087–1100. doi:10.1016/j.advwatres.2008.04.012
- Stedinger, J. R., *et al.*, 2008. Appraisal of the generalized likelihood uncertainty estimation (GLUE) method. *Water Resources Research*, 44, W00B06. doi:10.1029/2008WR006822
- Todini, E. and Mantovan, P., 2007. Comment on: “On undermining the science?” by Keith Beven. *Hydrological Processes*, 21 (12), 1633–1638. doi:10.1002/hyp.6670
- Vrugt, J. A. and Sadegh, M., 2013. Toward diagnostic model calibration and evaluation: approximate Bayesian computation. *Water Resources Research*, 49, 4335–4345. doi:10.1002/wrcr.20354
- Vrugt, J. A., *et al.*, 2008. Treatment of input uncertainty in hydrologic modeling: doing hydrology backward with Markov chain Monte Carlo simulation. *Water Resources Research*, 44 (12), W00B09. doi:10.1029/2007WR006720
- Westerberg, I., *et al.*, 2011a. Stage-discharge uncertainty derived with a non-stationary rating curve in the Choluteca River, Honduras. *Hydrological Processes*, 25, 603–613. doi:10.1002/hyp.7848
- Westerberg, I. K., *et al.*, 2011b. Calibration of hydrological models using flow-duration curves. *Hydrology and Earth System Sciences*, 15 (7), 2205–2227. doi:10.5194/hess-15-2205-2011
- Wicks, J. M., Adamson, M., and Horritt, M., 2008. Communicating uncertainty in flood maps—a practical approach. *Defra Flood and Coastal Management Conference*, Manchester.
- Wilby, R. L. and Dessai, S., 2010. Robust adaptation to climate change. *Weather*, 65 (7), 180–185.