

INFORMATION THEORY, UNCERTAINTY AND RISK FOR EVALUATING HYDROLOGIC FORECASTS

by

Steven Weijts⁽¹⁾ and Nick van de Giesen⁽²⁾

⁽¹⁾ Faculty of Civil and Environmental Engineering, Delft University of Technology, Delft, The Netherlands (s.v.weijts@tudelft.nl)

⁽²⁾ Faculty of Civil and Environmental Engineering, Delft University of Technology, Delft, The Netherlands
(n.c.vandegiesen@tudelft.nl)

ABSTRACT

Probabilistic forecasting is becoming increasingly popular in hydrology. Equally important are methods to evaluate such forecasts. There is still debate about which scores to use for this evaluation. In this paper we distinguish two scales for evaluation: information-uncertainty and utility-risk. We claim that the information-uncertainty scale is the most natural for forecast evaluation. We propose a Kullback-Leibler divergence as the appropriate measure for forecast quality and show that some commonly used scores are approximations of it. From a new decomposition of our information-theoretical score into uncertainty, correct information and wrong information, it follows directly that deterministic forecasts increase uncertainty to infinity, although they can still have value. Because information is uniquely defined, any measure that does not measure information, can be interpreted as some other form of utility. We claim that for calibration of models representing a hydrological system, information should be the objective, because it allows to learn more from the observations. Any other utility measure trains an implicit decision model, which inevitably results in a loss of information in the data and more risk of overfitting.

Keywords: information-theory, relative entropy, Kullback-Leibler divergence, forecast verification, utility, calibration

1 INTRODUCTION

Over the last decades, probabilistic forecasting has become increasingly important in the field of hydrology. Lacking enough information to completely eliminate uncertainty, probabilistic forecasts are intended to reduce uncertainty of the user about the future events and communicate the remaining uncertainty (Krzysztofowicz2001; Montanari and Brath2004). In hydrology, the development of methods for evaluating such forecasts, however, has not kept pace with the development of methods of generating them (Laio and Tamea2007). This is an important problem, given the fact that science is required to make testable predictions and therefore needs unambiguous methods for testing those predictions. Furthermore, the lack of methods for evaluation of hydrological forecasts may hinder acceptance of those forecasts by the public. In this paper we approach forecast evaluation from an information-theoretical point of view. By using a decomposition we developed recently, we provide some insights in what evaluation scores measure and what, in our opinion, they should measure.

1.1 What is a good forecast?

In general, the evaluation of forecasts can serve several purposes. Evaluation may serve to assign a level of trust in the forecast, to compare different forecasting systems or forecasters, to reward good forecasters, to diagnose problems in forecasting models, and to provide an objective function for calibration of the forecasting models. All these purposes for evaluation have in common that the measures should allow comparisons between forecasts or between series of forecasts. Assigning a level of trust only makes sense if there are also alternatives; rewarding a good forecaster has no use if there is no other forecaster or no other period of forecasts to compare to; diagnosing problems is not possible if there is no reference of what the quality should be; optimization works by continuously comparing different models or parameter sets.

For directly comparing two (series of) forecasts, preferences must be complete (a forecast must either be better, worse, or equally good than another one) and transitive (preferences can not form a loop like $A > B > C > A$), which are the same requirements that are applicable to probability. These two requirements

naturally lead to measures that take the form of a scalar real number. In contrast to this one-dimensional requirement, however, (Murphy1993) argued that it is possible to distinguish three different dimensions of forecast “goodness”:

- Consistency: correspondence between forecasts and judgements
- Quality: The correspondence between forecasts and observations
- Value: Incremental benefits of forecasts to users

Consistency requires that what the forecaster communicates, the forecasts, corresponds to his best judgements. Those judgements are internal to the forecaster and should be a rational distillation of all information available to him. Because a forecaster has only limited access to information and is not completely rational, *his* best judgements may not be *the* best judgements, but by definition he can never knowingly let *his* best estimate diverge from *the* best estimate, or it would not be his best estimate.

Quality is the dimension that is most important in pure science, as it concerns putting the predictions to the test by comparing forecasts with observations. It is important to note in this respect that an observation is also only an estimate of the truth and therefore does not fundamentally differ from a forecast. In fact, we are comparing one estimate of truth with another. The estimate that we regard as most trustworthy, usually the one that is made in hindsight, is called observation, the other estimate is the prediction or forecast. In meteorology, the evaluation of quality is called verification (latin: veritas = truthfulness). This term is somewhat misleading, because establishing that a model simulates the truth is impossible (Oreskes et al.1994).

Value is related to a decision problem attached to the forecast. It is therefore not only dependent on the forecasts and the observations, but also on who are using the forecasts. Hydrological forecasts may, for example, have significant value for reservoir operation, evacuation decisions, and agriculture. Good forecasts can lead to more hydropower, less flood damage, and, at the same time, less unnecessary pre-releases for flood protection. It could be attempted to express these benefits in monetary terms, but from a decision-theoretical point of view, it is better to use the more general term utility. This takes into account that not every unit of money necessarily has the same value. By definition, the utility of an uncertain outcome is equal to the expected utility of that outcome. In engineering, risk is defined as expected damage or disutility. Risk is therefore the opposite of utility. For adverse events, like floods, anticipation can reduce risk and the value of hydrological forecasts can thus be expressed as the reduction in risk they provide when used in decision making.

1.2 Problems with evaluation of hydrological forecasts

The current problem in defining a framework for evaluation of forecasts lies partly in that the distinction between the latter two dimensions is not always explicitly made. As certain purposes of evaluation require a one-dimensional measure of goodness, a choice between value and quality must be made and if the latter is chosen, an unambiguous quality measure must be defined that can not rely on user preferences. The hydrological and meteorological literature, however, offers a wide range of verification measures. Although the properties of these measures are well-studied, it is not always clear what is actually measured. Laio and Tamea (2007) give an overview of some commonly used measures in meteorology that could be applicable in hydrology.

What is missing from this overview, and also in two standard works about forecast verification (Wilks2005) and (Jolliffe and Stephenson2003), are measures for forecast evaluation based on information theory (Weijs et al., 2010). We argue that information-theoretical scores are measures for quality par excellence, for forecasts stated in terms of probability.

Except for probabilistic forecasts, two other types of forecasts were presented in *Table I* of Laio and Tamea (2007): deterministic forecasts and interval forecasts. We think that these types of forecasts can in principle not be evaluated unambiguously without reference to external assumptions relating to probability or utility. Instead of seeing this as a problem of the evaluation methods, this should be seen as a problem of the forecasts themselves. They do not fulfil the requirement of testable predictions. Moreover, deterministic forecasts are often not consistent with judgements, which, given that we know a model is an approximation, are better described in terms of probability.

Notwithstanding these problems with deterministic forecasts, they are still common in hydrology and are usually evaluated with measures like NSE and MSE, MAE. Actually, many of the methods for producing probabilistic forecasts make use of deterministic forecasts and their evaluation, for example Markov Chain Monte-Carlo based methods. Therefore, it is likely that there exists some reason that makes deterministic forecasts acceptable from a practical point of view. Also here the information-theoretical viewpoint could provide some new insights.

1.3 Outline

In this paper, we propose to use information theory as the central framework for forecast quality. By viewing the forecast evaluation problem from an information-theoretical perspective, we hope to shed some light on what is measured and what should be measured.

In section 2, we present an information-theoretical score for forecast quality and a decomposition that we recently found (Weijs et al, 2010). We also show that the components of a commonly used Brier score decomposition are second order approximations to our information components. Our information-theoretical divergence score can be interpreted as remaining uncertainty after receiving the forecast. In section 3 we analyse the seemingly paradoxical implication that deterministic forecasts increase the uncertainty to infinity and we offer two interpretations to resolve this paradox. In section 4, the question is addressed whether or not the utility a model provides for users should be considered in the calibration process. The conclusions are summarized in the last section, where we argue that issuing forecasts can best be considered a communication problem and that the information they provide is the most sensible measure for their evaluation.

2 INFORMATION-THEORETICAL EVALUATION OF FORECASTS

Information theory provides a number of measures relating uncertainty and information, within the framework of probability theory. Since forecasting can be seen as providing information to reduce uncertainty about future events, information-theory appears to be an appropriate framework to evaluate forecasts. As we showed in our recent paper, the measure Kullback-Leibler divergence, or relative entropy, can be used as a verification score and has a number of desirable properties. Starting from an analogy with the Brier score, we now introduce the divergence score and an insight-providing decomposition of it. For a more elaborate description and some other related discussions, see (Weijs et al., 2010)

2.1 Classical decomposition of the Brier Score

The Brier score was introduced by (Brier1950) as a verification score for probabilistic forecasts. It is still the most widely used score for evaluating probabilistic forecasts of binary events. A binary event has two possible outcomes, e.g. exceedence or non-exceedence of a certain critical water level in a river. A probabilistic forecast for one such a binary event at time t can be represented by a probability mass function (PMF), which in this case is a 2D vector, denoted by \mathbf{f}_t . The bold notation indicates a vector. When a probabilistic flow forecast indicates that there is 20% chance that the critical flow will be exceeded, for example, the forecast can be written as $\mathbf{f}_t = (1-f, f)^T = (0.8, 0.2)^T$, where the scalar f denotes the probability of exceedence. After the event is observed, the observation can also be written as a PMF, this time expressing the probabilities after the event has been observed. In case we assume perfect observations, and we observed exceedence of the critical level, the observation can be expressed as $\mathbf{o}_t = (1-o, o)^T = (0, 1)^T$. In this paper, we assume perfect observations to allow for the decompositions we introduce, but in general, perfect observations are no necessary assumption for the scores to be meaningful. Given the preceding definitions, the Brier score can now be defined as:

$$BS_t = 2(f_t - o_t)^2 = (\mathbf{f}_t - \mathbf{o}_t)^2 := (\mathbf{f}_t - \mathbf{o}_t)^T (\mathbf{f}_t - \mathbf{o}_t) \quad (1)$$

It must be noted that the Brier score is nowadays almost always defined as half this value (Ahrens and Walser2008). To make notation easier, we use the original definition of Brier (see eq. 1). For a series of forecasts, the Brier score is defined as the average of eq. 1 over all forecast instances. It can be interpreted as the mean squared error (MSE) in probabilities.

(Murphy1973) showed that the Brier score for such a series can be decomposed into three components: uncertainty, resolution and reliability:

$$BS = REL_{BS} - RES_{BS} + UNC_{BS} \quad (2)$$

$$BS = \frac{1}{N} \sum_{k=1}^K n_k (\mathbf{f}_k - \bar{\mathbf{o}}_k)^2 - \frac{1}{N} \sum_{k=1}^K n_k (\bar{\mathbf{o}}_k - \bar{\mathbf{o}})^2 + \bar{\mathbf{o}}(1 - \bar{\mathbf{o}}) \quad (3)$$

Where N is the total number and K the number of unique of forecasts issued, $\bar{\mathbf{o}} = \sum_{t=1}^N \mathbf{o} / N$ the climatological probability of occurrence of the event, n_k the number of forecasts within one category of unique forecasts and $\bar{\mathbf{o}}_k$ the observed frequency, given forecasts of probability \mathbf{f}_k .

The uncertainty term measures the inherent uncertainty in the climate. The uncertainty reaches a maximum for equiprobable outcomes and is zero if the outcome is always the same. The resolution and reliability terms in this decomposition can be seen as squared Euclidean distance measures between two probability distributions. The resolution term measures how much of the climatic uncertainty can be resolved by the forecasts. This is expressed in the average distance of the conditional distributions of the observations from the marginal distribution of the observations. The reliability measures the average squared distance between the forecast distributions and the corresponding conditional distributions of observations. A perfect reliability of zero (a more accurate term would be unreliability) is attained when for all forecast probabilities, the observed conditional frequency matches that probability. In this case the forecast is said to be perfectly calibrated.

2.2 Information-theoretical equivalents: Divergence score and decomposition

Information theory started with the paper of (Shannon1948), where he derived a measure of uncertainty (entropy) from three basic requirements for such a measure. The highly readable original paper is recommended for more background. The uncertainty of the climate using this definition is

$$H(\bar{\mathbf{o}}) = -\sum_{i=1}^n \left\{ \left[\bar{\mathbf{o}} \right]_i \log \left[\bar{\mathbf{o}} \right]_i \right\} \quad (4)$$

The logarithm has base 2, yielding the measure H in the unit bits. A related measure is relative entropy, also known as Kullback-Leibler divergence. This is a measure of the extra amount of uncertainty if one distribution is assumed, while the true distribution is different, this is the divergence from the true to the other distribution. In contrast to a distance like the Brier score, Kullback-Leibler divergence is not symmetric. The divergence depends on which of the two distributions is considered the true one.

We define the divergence score as the divergence *from* the observation PMF *to* the forecast PMF:

$$DS_t = D_{KL}(\mathbf{o}_t \parallel \mathbf{f}_t) = \sum_{i=1}^n [\mathbf{o}_t]_i \log \left(\frac{[\mathbf{o}_t]_i}{[\mathbf{f}_t]_i} \right) \quad (5)$$

Where n is the number of possible outcomes (2 in the binary case). For a series of N forecasts and corresponding observations, the divergence score is

$$DS = \frac{1}{N} \sum_{t=1}^N D_{KL}(\mathbf{o}_t \parallel \mathbf{f}_t) \quad (6)$$

When replacing all quadratic distances in the Brier score decomposition by the appropriate divergences and replacing the uncertainty component by the information-theoretical definition of uncertainty, entropy, we obtain (see *Table I*):

$$DS = \frac{1}{N} \sum_{t=1}^N D_{KL}(\mathbf{o}_t \parallel \mathbf{f}_t) = \frac{1}{N} \sum_{k=1}^K n_k D_{KL}(\bar{\mathbf{o}}_k \parallel \mathbf{f}_k) - \frac{1}{N} \sum_{k=1}^K n_k D_{KL}(\bar{\mathbf{o}}_k \parallel \bar{\mathbf{o}}) + H(\bar{\mathbf{o}}) \quad (7)$$

In the appendix of (Weijs et al., 2010) it was shown that this equation holds, and thus we have obtained an information-theoretical equivalent of the Brier score and its decomposition.

2.3 Relation between the divergence and the Brier score

The components of the Brier score are second order approximations of the components of the divergence score (fig. 1). The uncertainty has the same location of maximum and zero points. When scaled with its maximum value, the similarity becomes visible (see left fig. 1). The resolution (right fig. 1), can reach a maximum equal to the uncertainty term. When scaled with the uncertainty, again a similarity between the shapes of the the resolution components is visible. The reliability term, however, exhibits significant differences in the extremes. While the reliability term of the Brier score is bounded, the analogous term in the divergence score can reach infinity. This happens when an outcome happens that was given zero probability in the forecast.

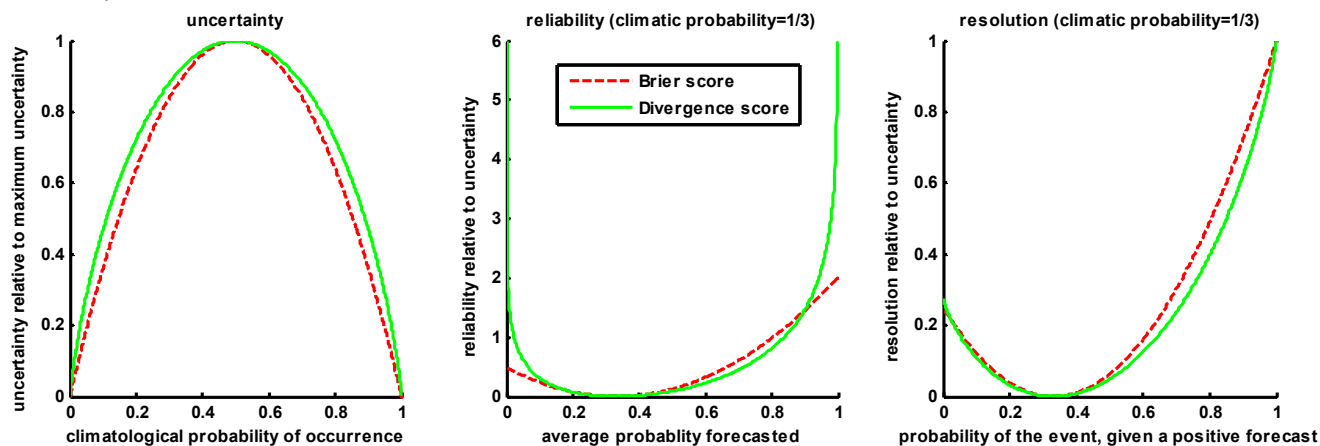


Figure 1 – The Brier score components and second order approximations of the divergence score components.

Table 1 – The expressions for the two decompositions compared

| | UNC | REL | RES |
|-------------------------|----------------------|---|--|
| Brier Score | $\bar{o}(1-\bar{o})$ | $\frac{1}{N} \sum_{k=1}^K n_k (\mathbf{f}_k - \bar{o}_k)^2$ | $\frac{1}{N} \sum_{k=1}^K n_k (\bar{o}_k - \bar{o})^2$ |
| Divergence Score | $H(\bar{o})$ | $\frac{1}{N} \sum_{k=1}^K n_k D_{KL}(\bar{o}_k \parallel \mathbf{f}_k)$ | $E_k \{D_{KL}(\bar{o}_k \parallel \bar{o})\}$ |

2.4 Interpretations of divergence score and decomposition

The new information-theoretical equivalents of the components of the Brier score allow some additional interpretations. One of the interpretations of measures in information-theory starts from a definition of surprise. Surprise is something we feel when something unexpected happens. The lower the probability we assume something to have, the more surprised we are when observing it. Rain in a desert is surprising, rain in the Netherlands is less surprising and rain on the moon is a miracle yielding almost unbounded surprise. When, the surprise of observing outcome x is defined as $S_x = \log(1/P(x))$, surprise can be measured in bits like information and uncertainty. Observing something that was a certain fact yields no surprise, heads on a fair coin yield one bit of surprise and observing a 1/1000 year flood in some year yields a surprise of approximately 10 bits. The entropy-measure for uncertainty can now be interpreted as the expected surprise about the truth: $H(X) = E_X(S_x)$

In general, uncertainty can now be interpreted as expected surprise about the true outcome. The fact that different expectations can be calculated according to different subjective probability distributions, reflects that uncertainty can be both something objective and subjective. The uncertainty a person thinks to have is the entropy of his subjective probability distribution. Kullback-Leibler divergence can be seen as the additional uncertainty one person estimates the other person to have compared to his own:

$$D_{KL}(P(X) \parallel Q(X)) = E_{P(X)}(S_{Q(X)} - S_{P(X)})$$

Because forecast verification is done in hindsight, the observation that is made can be used as a viewpoint to estimate the uncertainty. The additional uncertainty (surprise about the truth), estimated from the viewpoint of the observation is the best available estimate of the remaining uncertainty about the truth of the person having the forecast. Assuming perfect observations, the divergence score measures remaining uncertainty about the truth. In fig. 2 it is shown how the components of the divergence score relate to this remaining uncertainty at different levels of informedness. Interpreting the figure, resolution can be seen as the correct information, that can be subtracted from the climatological uncertainty or missing information. The reliability term is added to the remaining uncertainty and represents the wrong information due to biased probability estimates. The wrong information can be reduced by calibration. It should be noted that the decomposition is only meaningful when enough data is available to properly calculate all conditionals (Weijs, 2010).

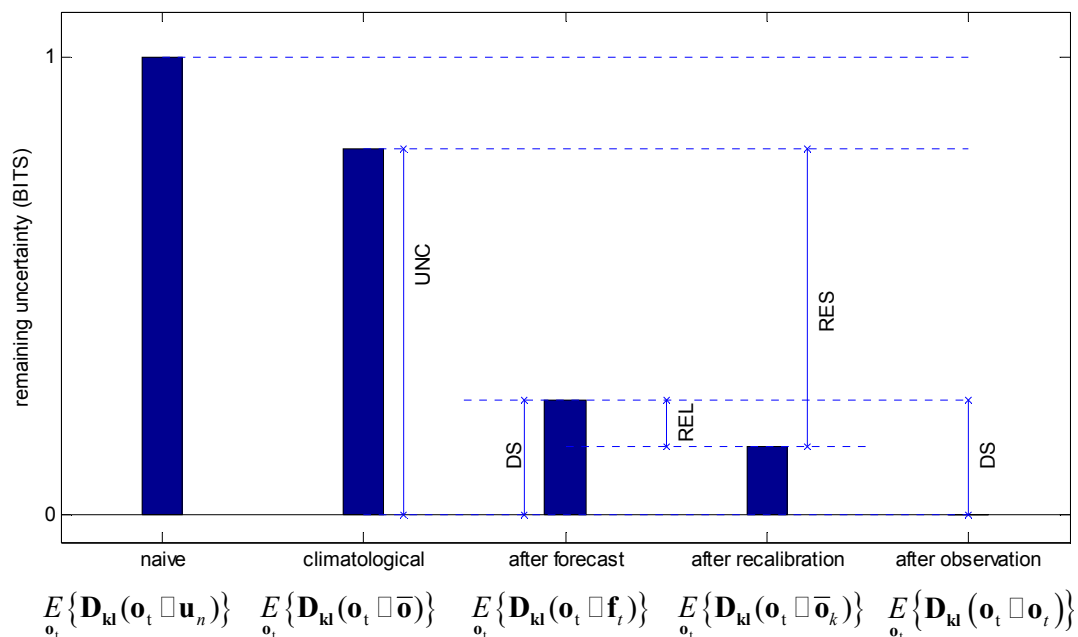


Figure 2 – The remaining uncertainty for different distributions in the forecasting proces.

3 DETERMINISTIC FORECASTS ARE INCONSISTENT

Can a forecaster be completely sure about something that in the end does not happen and still get credit for his forecast? This does not appear natural, but it often turns out to happen in practice. For example, a deterministic flow forecast of 200 m³/s is considered quite good, when 210 m³/s is observed. Apparently, it is already expected that some error will occur and a forecast that is 10 m³/s off is considered to be not that bad. Hydrological models are per definition simplifications of reality. Often, they describe relations between macrostates, like averaged rainfall, mass of water in the groundwater reservoir, and flow through a river cross-section. Just like problems in statistical thermodynamics, having limited information about what really goes on inside a hydrological system on a microscopical level, our forecasts can never be perfect (Weijs2009). What can be said about the real world on the basis of a model is therefore inherently erroneous to some extent, or should be stated in terms of probabilities.

How then, should deterministic forecasts be evaluated? Literally taken, a deterministic (point value) forecast states: “the outcome is x”. Implicitly, such a forecast asks to be evaluated from a black and white view: the forecast is either wrong or right. The divergence score also reflects this. If the forecast was right, the perfect score of 0 is attained, if the forecast was wrong, however, a penalty of infinity is given. If one such a forecast is given, the forecaster can look for another career, because even a future series of perfect forecasts can not average out the infinite penalty. The decomposition shows that the reliability component is responsible. Although the deterministic forecasts usually contain information about the observed outcomes, given that the resolution is positive and removes some of the uncertainty, this is completely annihilated by the reliability term. The discrepancy between the information (reduction of uncertainty) that the forecasts contain and the information conveyed by the messages that constitute the forecasts is so large that the

expected surprise about the truth of a person taking the forecast at face value goes to infinity. The fact that deterministic forecasts are still used in society (and unfortunately sometimes even preferred), while they explode uncertainty to infinity, seems to present a paradox. We propose two possible interpretations that offer a solution to this paradox.

3.1 Deterministic forecasts are implicitly probabilistic (information interpretation)

Fortunately, in reality, almost no person using deterministic forecasts takes them at face value. In fact, the forecast is implicitly recalibrated by the user, reducing the reliability term for the internal probability estimates the user bases his actions on. The user can do the recalibration based on previous experience with the forecasts and common sense. The user of the forecast can think “if the forecaster says the water level will be 10 cm under the embankment, he implicitly also forecasts a little that overtopping will occur”. Note that the example of Grand Forks in (Krzysztofowicz2001) shows that not all users do this. Mathematically this recalibration is equivalent to also attaching some probability to overtopping. However, it is not the task of a user to guess what the forecaster wanted to say. Consistency requires that the forecaster communicates his judgements to the user (Murphy1993).

The forecaster may also present the forecast as being an expected value or mean. This suggests an underlying probabilistic forecast. However, when taking the information-theoretical viewpoint, communicating an expected value means nothing without additional statements regarding the probability distribution. The principle of maximum entropy (PME) (Jaynes1957) states that when making inferences based on incomplete information, the best estimate for the probabilities is the distribution that is consistent with all information, but maximizes uncertainty. In this way, the uncertainty is reduced exactly by the amount the information permits, but not more. Maximizing entropy with known mean and variance, gives a Gaussian distribution, maximizing uncertainty about the velocities of gas molecules with known total kinetic energy gives the Boltzmann distribution. When PME is applied to expected value forecasts, however, the maximum entropy forecast distribution that is consistent with the information given by the forecaster is uniform between minus and plus infinity. It is the complete opposite end of the spectrum compared to the previous literal interpretation of the deterministic forecast: from claiming total certainty to claiming total uncertainty.

In the case of streamflow forecasts, the user can still get a less nonsensical forecast distribution by combining the information in the forecast with the common sense notion that streamflows in rivers are nonnegative. This extra constraint turns the PME forecast distribution for a known expected value into an exponential distribution.

When a series of deterministic forecasts is evaluated, it is possible to look at the joint distribution of forecasts and observations. The conditional distributions of observations for each forecast value can then be seen as probabilistic forecast distributions. It is important to note however, that the probabilistic part of such a forecast is derived from data. When such forecasts are evaluated, the predictive performance is judged on the basis of an error model, that is derived from the same data that is used for its evaluation. For example, if a deterministic model is evaluated with mean squared error, the evaluation can also be interpreted as evaluating something proportional to the divergence score (remaining uncertainty), if the deterministic forecasts are assumed to implicitly convey probabilistic forecasts, being normal distributions with the mean at the deterministic forecast value. The second parameter (variance) of this distribution, however, is implicitly assumed to have the best fit to the errors. While this approach may under some conditions be acceptable for calibration to train the error model, for evaluation of forecasts it is unacceptable, because it uses the data against which it is evaluated. A correct approach would be to explicitly formulate and train an error model in the calibration, and use that model to make probabilistic predictions for the evaluation period, that can subsequently be evaluated with the divergence score. The error models are not restricted to Gaussian distributions, but can take more flexible forms. Such an approach is taken in Schoups et al. (2010).

As a last consideration, we want to stress that even if an error model is properly formulated and added to the deterministic “physical” part, the resulting model still represents a false dichotomy between true behaviour of the system and the error, as was argued by (Koutsoyiannis2009). A more consistent approach would be to explicitly make the probabilistic part of the model an integrated part of the physical reality it is supposed to simplify. Such approaches can lie in studying the time-evolution of chaotic systems (Koutsoyiannis2009) or

in applying the principle of maximum entropy in combination with macroscopic constraints, as suggested by Weijs2009).

Concluding, from the information-theoretical viewpoint several reasons come to light why deterministic forecasts should in fact be considered to be probabilistic. The problem with these forecasts is that they leave too much of the probabilistic interpretation to the user. It might be considered ironic that the users who are claimed to not be able to handle probabilistic forecasts and are for that reason provided with deterministic forecasts are the ones who have to rely most on their ability to subconsciously make probability estimates based on the limited information in the forecast.

3.2 Deterministic forecasts still can have value for decisions (utility interpretation)

A second, independent interpretation of deterministic forecasts that justifies their existence is their usefulness, even to users who do not make subconscious probability estimates. Even though a reservoir operator might be infinitely surprised if he has taken a deterministic inflow forecast of 200 m³/s at face-value and he finds out the inflow was 250 m³/s, his loss is not infinite. The operator might spill some water, but not all is lost.

The difference between surprise and loss is due to the fact that most decision problems are not equal to placing stakes in a series of horse races. Such a horse race is the classical example where information can be directly related to utility, see (Kelly1956) and (Cover and Thomas2006) for more explanation. In such a horse race, all money not bet on the winning horse is lost, so the only important probability is the one attached to the winning horse, that determines the stakes the gambler should put on that horse (Kelly showed that the stakes on each horse should be proportional to the estimated winning probabilities). In contrast, for decision problems like reservoir operation, optimally preparing for 200 m³/s automatically implies also preparing for 250 m³/s to some extent. This makes the loss function non-local (locality is discussed in Sec 4.1).

Another difference with the horse race is that the total amount of value at stake in hydrological decision making usually does not depend on the previous results, while the results for the horse race assume that the gambler invests all his previously accumulated capital in the bets. The gambler therefore wants to maximize the product of rates of return over the whole series of bets, while for a reservoir operator, each period offers a new opportunity to gain something from the water, even though he spilled all his water in the previous month. This is comparable with a gambler whose wife allows him to bet a fixed amount of money each week (Kelly1956) and then spends it all in the bar on the same evening without possibility of reinvesting in the next bet. Assuming a utility linear in the beer he buys with the winnings, the best decision is to bet all money on the one horse with the best expected return. Again, one loss is not fatal for the whole series of bets. He just hopes for better luck next week. The evaluation of the value of deterministic forecasts is therefore not as black and white as evaluation of the information they contain.

The evaluation of deterministic forecasts in this interpretation is thus connected to a decision problem. Decisions can be taken as if the forecasts are really certain, and still be of value. The loss functions for evaluating forecasts can be seen as functions that somehow map the discrepancy between forecast value and observed value to a disutility of the decision based on the wrong forecast, compared to a perfect forecast. In the utility interpretation, evaluating deterministic forecasts with mean squared error implicitly defines a decision process in which the disutility is a quadratic function of the distance between forecast and observation. In that case, forecasts that have the smallest MSE have most utility or value for the user.

4 INFORMATION VERSUS UTILITY AS CALIBRATION OBJECTIVE

Value-based forecast evaluation is inevitably connected to a particular user with a decision problem and therefore cannot be done without explicit consideration of the user base of forecasts. Moreover, an obvious question that arises is whether it is desirable to base the evaluation on the value to a particular user or group of users. In that case, the evaluation becomes an evaluation of decisions rather than of the forecasts themselves or of the hydrological model that produced them. This difference is particularly important if the results of the evaluation are used in a learning or calibration process.

4.1 Locality and philosophy of science

Locality is a property of scores for probabilistic forecasts. A score is said to be local if the score only depends on the probability assigned to the event that occurred, and does not depend on how the probability

is spread out over the values that did not occur. In contrast to this, non-local scores do depend on how that probability is spread out. Usually they are required to be sensitive to distance, which means that probability attached to values far from the observed value is punished more heavily than forecast probability that was assigned to values close to the observation. This concept of distance only plays a role in forecasts of continuous and ordinal discrete predictands. For both these types of predictands, an extension of the Brier score exists: the Ranked Probability Score (RPS) and the continuous RPS (CRPS) (See Laio and Tamea, 2007 for description and references). Both these scores are non-local, while the divergence score is local.

For most decision problems, expected utility is a non-local score: a reservoir operator that attached most probability to values far from the true inflow is worse off than one that used a forecast with most probability close to the true value, even if the probability (density) attached to the true value was the same. Therefore, non-local scores are sometimes considered to have more intuitive appeal than local scores.

There is, however, a serious philosophical problem with non-local scores if used in a learning process. It is a violation of scientific logic if the score that is intended to evaluate the quality of forecasts depends on what is stated about things that are not observed. In an extreme case, two series that forecast the same probabilities for all the events that were observed, can obtain different scores based on differences in the probabilities assigned to unobserved events (Benedetti2010). A similar argument in the context of experimental design was made by (Bernardo1979) If these non-local scores are used as objectives in calibration or inference (see for example Gneiting et al.2005), things are thus inferred from non-observed outcomes, i.e. information that is not there.

4.2 Utility as a data filter

The use of utility in calibration can, next to using non-existing information, also lead to learning only from part of the information that is in the observations. In that sense, the decision problem that specifies the utility acts like a filter on the information. For example, when a binary evacuation decision is coupled to a conceptual rainfall-runoff model for flood forecasting, the calibration towards maximum utility of the system will train the hydrological model to optimally distinguish flood-evacuation events. This implies that in the training, all that the hydrological model sees from the continuous observed discharges is a binary signal: flood or no flood. This constitutes at most one bit of information per observation (in the unlikely case that a 50% of the forecasts leads to an evacuation decision), while the original signal contained far more information. The hydrological model will therefore have far less information to learn from.

Given the fact that there is a balance between the available information for calibration and the complexity that a model is allowed to have, hydrological models that are trained on this kind of utility functions are likely to become overly complex relative to the data (see Schoups, 2008). They will surely achieve better utility results on the calibration data (because there is less information to fit), but are likely to perform worse on an independent validation dataset. The model that has been trained with maximum information as an objective is likely to yield better results for the validation set, even in terms of utility. Because it has the unfiltered information from the observations to learn from, it is less prone to overfitting: the complexity of a conceptual hydrological model is better warranted by the full information. Training for optimal classification of flood events would benefit from more parsimonious data-driven models that make a mapping directly from predictors to decisions, but this complicates the use of prior information on the workings of the hydrological system.

The information-theoretical logarithmic scoring rules are the only scoring rules that are both local and proper (proofs can be found in Bernardo, 1979 and Benedetti, 2010). Where propriety is the requirement that the scoring rule can only be optimized when the forecaster does not lie. All utility functions that are not linear functions of information violate either locality or propriety, which makes them doubtful objectives for calibration.

5 CONCLUSIONS

The difficulties and debate about the evaluation of forecasts can be significantly clarified using results from information theory. When information is seen as a measurable quantity, like energy, a sort of “information intuition” develops, similar to the “energy intuition” that is used to detect logical flaws in claims for *perpetuum mobiles*. Science is required to make testable predictions. Forecasts should therefore be stated in terms that make it clear how to evaluate them. Deterministic and interval forecasts fail this criterion.

Probabilistic forecasts can be evaluated using information theory. The decomposition of the divergence score that we presented can provide additional insight in the interaction between uncertainty, correct information and wrong information.

Starting from the observation that deterministic forecasts are still commonly used and evaluated, but are worthless from an information-theoretical viewpoint, we draw the conclusion that these forecasts are either implicitly probabilistic or should be viewed in connection to a decision problem. In both interpretations, the evaluation depends on external information that is not provided in the forecast. Deterministic forecasts leave too much interpretation to the user, if seen as implicit probabilistic forecasts or make too many assumptions on the user if they are evaluated using another utility measure.

On the one hand, forecasting can be seen as a communication problem in which uncertainty about the outcome of a random event is reduced by delivering an informative message to a user. On the other hand, forecasting can be seen as an addition of value to a decision problem. Any measure that is not information only becomes meaningful when it is interpreted in terms of utilities. When addressing forecast value, it is important to see that in fact we are evaluating decisions based on forecasts and not the correspondence between the observations and the forecasts themselves.

This is especially important in calibration, where a model has to learn from observations. When calibration objectives are used that are not information-measures, the model either learns from information that is not there or uses only part of the information in the observations, or both. Because the amount of available information is related to optimal model complexity, hydrological models trained for user specific utilities are more prone to overfitting, which might lead to worse results in an independent validation test.

6 REFERENCES

- Ahrens, B. and Walser, A. (2008). Information-based skill scores for probabilistic forecasts. *Monthly Weather Review*, 136(1):352–363.
- Benedetti, R. (2010). Scoring Rules for Forecast Verification. *Monthly Weather Review*, 138:203.
- Bernardo, J. (1979). Expected information as expected utility. *The Annals of Statistics*, 7(3):686–690.
- Brier, G. W. (1950). Verification of forecasts expressed in terms of probability. *Monthly weather review*, 78(1):1–3.
- Cover, T. and Thomas, J. (2006). *Elements of information theory*. Wiley-Interscience.
- Gneiting, T., Raftery, A., Westveld, A., and Goldman, T. (2005). Calibrated probabilistic forecasting using ensemble model output statistics and minimum CRPS estimation. *Monthly Weather Review*, 133(5):1098–1118.
- Jaynes, E. T. (1957). Information theory and statistical mechanics. *Physical Review*, 106(4):620–630.
- Jolliffe, I. T. and Stephenson, D. B. (2003). *Forecast verification: a practitioner's guide in atmospheric science*. Wiley.
- Kelly, J. (1956). A new interpretation of information rate. *Information Theory, IEEE Transactions on*, 2(3):185–189.
- Koutsoyiannis, D. (2009). HESS Opinions: A random walk on water. *Hydrology and Earth System Sciences Discussions*, 6:6611–6658.
- Krzysztofowicz, R. (2001). The case for probabilistic forecasting in hydrology. *Journal of Hydrology*, 249(1):2–9.
- Laio, F. and Tamea, S. (2007). Verification tools for probabilistic forecasts of continuous hydrological variables. *Hydrology and Earth System Sciences*, 11(4):1267–1277.
- Montanari, A. and Brath, A. (2004). A stochastic approach for assessing the uncertainty of rainfall-runoff simulations. *Water Resources Research*, 40:1106.
- Murphy, A. H. (1973). A new vector partition of the probability score. *Journal of Applied Meteorology*, 12(4):595–600.
- Murphy, A. H. (1993). What is a good forecast?: An essay on the nature of goodness in weather forecasting. *Weather and Forecasting*, 8(2):281–293.
- Oreskes, N., Shrader-Frechette, K., and Belitz, K. (1994). Verification, validation, and confirmation of numerical models in the earth sciences. *Science*, 263(5147):641.
- Schoups, G., N. C. van de Giesen, and H. H. G. Savenije (2008), Model complexity control for hydrologic prediction, *Water Resources Research*, 44(1)
- Shannon, C. E. (1948). A mathematical theory of communication. *Bell System Technical J.*, 27(3):379–423.
- Weijs, S. (2009). Interactive comment on “HESS Opinions : A random walk on water” by D. Koutsoyiannis. *Hydrology and Earth System Sciences Discussions*, 6:2733–2745.
- Weijs, S., Van Nooijen, R., and Van de Giesen, N. (2010). Kullback–Leibler divergence as a forecast skill score with classical reliability–resolution–uncertainty-decomposition. *Monthly Weather Review*, accepted.

Wilks, D. S. (2005). *Statistical methods in the atmospheric sciences*. 2nd edition.