

Predictive Uncertainty in Environmental Modelling

Gavin C. Cawley

School of Computing Sciences

University of East Anglia

Norwich NR4 7TJ

United Kingdom

E-mail: gcc@cmp.uea.ac.uk

Malcolm R. Haylock

Climatic Research Unit

University of East Anglia

Norwich NR4 7TJ

United Kingdom

E-mail: M.Haylock@uea.ac.uk

Stephen R. Dorling

School of Environmental Sciences

University of East Anglia

Norwich NR4 7TJ

United Kingdom

E-mail: S.Dorling@uea.ac.uk

Abstract—Artificial neural networks have proved an attractive approach to non-linear regression problems arising in environmental modelling, such as statistical downscaling, short-term forecasting of atmospheric pollutant concentrations and rainfall run-off modelling. However, environmental datasets are frequently very noisy and characterised by a noise process that may be heteroscedastic (having input dependent variance) and/or non-Gaussian. The aim of this paper is to review an existing methodology for estimating predictive uncertainty in such situations, and more importantly illustrate how a model of the predictive distribution may be exploited in assessing the possible impacts of climate change and to improve current decision making processes. The results of the WCCI-2006 predictive uncertainty in environmental modelling challenge are also reviewed and some areas suggested where further research may provide significant benefits.

I. INTRODUCTION

Neural networks have been shown to provide a simple and flexible approach to non-linear regression problems arising in the environmental sciences. Some recent applications include statistical downscaling [13], water level-discharge modelling [1], river stage forecasting [8] and air quality forecasting [24]. The presence of special sessions devoted to environmental sciences and climate modelling at IJCNN-2005 and IJCNN-2006 provides further evidence of the importance of this field of research. Environmental modelling problems are typically very noisy and often characterised by a noise process that is heteroscedastic (i.e. the variance of the noise process is input-dependent) and may also be non-Gaussian. Conventional neural network regression techniques aim to estimate the conditional mean of the target data, via minimisation of a sum-of-squares error function. The aim of this paper is to demonstrate that practical benefits can be accrued by attempting to model the entire distribution of the noise contaminating the data in addition to the conditional mean. For example, we may estimate the conditional variance of the noise process, which may be achieved by training a second regression network to predict the squared residuals of the first. The combined model provides a Gaussian *predictive distribution* indicating the relative plausibility of different values for the target function. The provision of a predictive distribution, instead of only the conditional mean, can be exploited in a number of ways:

- The predictive distribution implies a plausible interval

(a.k.a. “error bars”) on all predictions, which in turn provide a valuable indicator of the reliability of the model.

- An estimate of the predictive distribution allows the estimation of the true *risk*, i.e. we may integrate the loss associated with all plausible outcomes, weighted by the probability of their occurrence.
- Where a neural network is used as one component within a much larger model, the uncertainties associated with the inputs and outputs of each component, may be propagated through the model (e.g. via a Monte-Carlo simulation) so that all sources of uncertainty can be integrated over to obtain a moderated prediction.
- Often we are interested in extreme events, especially the exceedance of some arbitrary threshold. For instance a statistical downscaling model might be used in estimating the impacts of future climate on the risk of flooding in a particular catchment, by modelling the linkage between large scale circulation and local precipitation. By their very nature, extreme events are not modelled well by an estimate of conditional mean of the data. However, given a full predictive distribution, we may at least estimate the *probability* of an extreme event by integrating the upper tail of the predictive distribution, even if the estimate of the conditional mean never exceeds the threshold.

Modelling predictive uncertainty in environmental data is also interesting from a machine learning perspective as the noise processes involved are often non-Gaussian and/or heteroscedastic, and so “off-the-shelf” solutions may not be entirely satisfactory, and thus there is significant scope for further research.

The remainder of this paper is structured as follows: Section II describes a simple methodology for estimating the predictive distribution based on methods developed by Peter Williams [29–32]. Section III demonstrates that an estimate of the predictive distribution can be exploited to provide practical benefits for the end-user, via an illustrative (if a little contrived) example based on the estimation of insurance losses associated with flood hazards. The results of the WCCI-2006 Predictive Uncertainty in Environmental Modelling Competition, which aimed to stimulate research

in this area, are presented in Section IV. Section V discusses some areas where further research may provide significant benefits. Finally, the work is summarised and conclusions drawn in Section VI.

II. MODELLING PREDICTIVE UNCERTAINTY WITH NEURAL NETWORKS

In this section, we outline a neural network approach to modelling predictive uncertainty in environmental applications, based on a sound Bayesian methodology developed by Williams [29–32]. For this study, we adopt the familiar Multi-Layer Perceptron network architecture (see e.g. Bishop [3]). The optimal model parameters, \mathbf{w} , are determined by gradient descent optimisation of an appropriate error function, $E_{\mathcal{D}}$, over a set of training examples, $\mathcal{D} = \{(\mathbf{x}_i, t_i)\}_{i=1}^N$, $\mathbf{x}_i \in \mathcal{X} \subset \mathbb{R}^d$, $t_i \in \mathbb{R}$, where \mathbf{x}_i is the vector of explanatory variables and t_i is the desired output for the i^{th} training pattern. The error metric most commonly encountered in non-linear regression is the sum-of-squares error, given by

$$E_{\mathcal{D}} = \frac{1}{2} \sum_{i=1}^N (y_i - t_i)^2, \quad (1)$$

where y_i is the output of the network for the i^{th} training pattern. In order to avoid over-fitting to the training data, however, it is common to adopt a regularised [27] error function, adding a term $E_{\mathcal{W}}$ penalising overly-complex models, i.e.

$$M = \alpha E_{\mathcal{W}} + \beta E_{\mathcal{D}}, \quad (2)$$

where α and β are regularisation parameters controlling the bias-variance trade-off [9]. Minimising a regularised error function of this nature is equivalent to the Bayesian approach which seeks to maximise the posterior density of the weights (e.g. [18, 20]), given by

$$P(\mathbf{w} \mid \mathcal{D}) \propto P(\mathcal{D} \mid \mathbf{w})P(\mathbf{w}),$$

where $P(\mathcal{D} \mid \mathbf{w})$ is the likelihood of the data and $P(\mathbf{w})$ is a prior distribution over \mathbf{w} . The form of the functions $E_{\mathcal{D}}$ and $E_{\mathcal{W}}$ correspond to distributional assumptions regarding the data likelihood and prior distribution over network parameters respectively. The usual sum-of-squares metric (1), corresponds to a Gaussian likelihood,

$$P(\mathcal{D} \mid \mathbf{w}) = \frac{1}{\sqrt{2\pi\beta^{-1}}} \exp\left\{-\frac{[t_i - y(\mathbf{x}_i)]^2}{2\beta^{-1}}\right\}$$

with fixed variance $\sigma^2 = 1/\beta$. Here, we adopt the Laplace prior propounded by Williams [30], which corresponds to a L_1 norm regularisation term,

$$E_{\mathcal{W}} = \sum_{i=1}^W |w_i|. \quad \iff \quad P(\mathbf{w}) = \frac{1}{2\beta} \exp\left\{-\frac{|w|}{\beta}\right\}$$

where W is the number of model parameters. An interesting feature of the Laplace regulariser is that it leads to the

automatic pruning of redundant model parameters. From 2, at a minimum of M we have

$$\left|\frac{\partial E_y}{\partial w_i}\right| = \frac{\alpha}{\beta} \quad w_i > 0, \quad \left|\frac{\partial E_y}{\partial w_i}\right| < \frac{\alpha}{\beta} \quad w_i = 0.$$

As a result, any weight not obtaining a data misfit sensitivity of α/β is set exactly to zero and can be pruned from the network.

A. Eliminating Regularisation Parameters

The hyper-parameters α and β can be estimated by maximising the evidence [18] or alternatively may be integrated out analytically [5, 30]. Here we take the latter approach; the posterior distribution of the parameters is given by

$$p(\mathbf{w}) = \int p(\mathbf{w}|\alpha)p(\alpha)d\alpha. \quad (3)$$

Assuming the Laplace prior, the prior distribution over the weights of the network, conditioned on the regularisation parameter α , is given by,

$$p(\mathbf{w}|\alpha) = Z_{\mathcal{W}}(\alpha)^{-1} \exp\{-\alpha E_{\mathcal{W}}\} \quad (4)$$

where the necessary normalising constant is given by

$$Z_{\mathcal{W}}(\alpha) = \left(\frac{2}{\alpha}\right)^W. \quad (5)$$

Substituting equations 4 and 5 into equation 3, adopting the (improper) uninformative Jeffreys prior, $p(\alpha) = 1/\alpha$ [16], and noting that α is strictly positive,

$$p(\mathbf{w}) = \int_0^\infty 2^{-W} \alpha^{W-1} \exp\{-\alpha E_{\mathcal{W}}\} d\alpha.$$

Using the Gamma integral, $\int_0^\infty x^{\nu-1} e^{-\mu x} dx = \frac{\Gamma(\nu)}{\mu^\nu}$ (Gradshteyn and Ryzhik [12], equation 3.384), we obtain

$$p(\mathbf{w}) = \frac{\Gamma(W)}{(2E_{\mathcal{W}})^W}.$$

Taking the negative logarithm and omitting irrelevant constant terms,

$$-\log p(\mathbf{w}) = W \log E_{\mathcal{W}}. \quad (6)$$

Applying a similar treatment to the data misfit term (assuming a sum-of-squares error), we have

$$L = \frac{1}{2} N \log E_{\mathcal{D}} + W \log E_{\mathcal{W}}.$$

For a network with more than one output unit, it is sensible to assume that each output has a different noise process (and therefore a different optimal value for β). It is also sensible to assign hidden layer weights and weights associated with each output unit to different regularisation classes so they are regularised separately. This leads to the training criterion used in this study:

$$L = \frac{N}{2} \sum_{i=1}^O \log E_{\mathcal{D}}^i + \sum_{j=1}^C W_j \log E_{\mathcal{W}}^j,$$

where O is the number of output units, C is the number of regularisation classes (groups of weights sharing the same regularisation parameter) and W_j is the number of non-zero weights in the j^{th} class. Note that bias parameters are not normally regularised. This approach provides a sound basic approach to non-linear regression using multi-layer perceptron networks, with Bayesian regularisation to prevent over-fitting *and* automatic selection of an appropriate network architecture as a result of the Laplace prior. As the regularisation parameters are integrated out analytically, the user need only select the initial number of hidden layer units, and more importantly an appropriate data misfit term that represents any available prior knowledge regarding the form of the noise process contaminating the data.

B. Choice of Data Misfit Term

In this paper, we are concerned with modelling predictive uncertainty, and so rather than simply estimating the conditional mean of the target data, we seek to construct a model such that the output specifies the entire predictive distribution. A sensible first step in solving an inference problem is to select an appropriate likelihood function to describe the statistical properties of the target data (c.f. [19]). The training criterion for the neural network should then be based on the negative logarithm of a parametric likelihood function, that incorporates any distributional assumptions regarding the noise process suggested by our prior knowledge of the data. In order to obtain a predictive distribution, we simply construct a network with one output for each of the parameters of this likelihood.

The most basic likelihood used in this study, assumes a heteroscedastic (input dependent variance) Gaussian noise process, i.e.

$$E_{\mathcal{D}} = \sum_{i=1}^{\ell} \left\{ \log \sigma(\mathbf{x}_i) + \frac{[\mu(\mathbf{x}_i) - t_i]^2}{2\sigma^2(\mathbf{x}_i)} \right\}. \quad (7)$$

Note the multi-layer perceptron network now has two output units, one giving the conditional mean of the target distribution, $\mu(\mathbf{x})$, as before, and an additional unit giving the conditional standard deviation, $\sigma(\mathbf{x})$. A linear activation function is used in the output unit corresponding to $\mu(\mathbf{x})$, and an exponential activation function for the unit corresponding to $\sigma(\mathbf{x})$, to enforce strictly positive estimates of conditional variance. This approach provides two advantages: Firstly the estimates of conditional variance provide error bars, indicating the uncertainty of model predictions [21, 22, 31]. Secondly the output of the model now completely specifies the target distribution, so the regularisation parameter β is no longer necessary. This data-misfit term is appropriate for regression on temperature data, where a Gaussian noise process is intuitively reasonable, but where the variability in temperature as well as the expected temperature may depend on, for example, the time of year.

The concentration of atmospheric pollutants provides an example of a type of data where a more complex likelihood

may be appropriate. Clearly a pollutant concentration cannot be negative, and the uncertainty in predictions is likely to be skewed upward. A common ploy would be to implement a log-normal likelihood, by simply taking the logarithm of the target data and employing the data misfit given in (7).

Modelling frontal precipitation data requires a more sophisticated statistical model, and is often modelled using a Gamma distribution [26] or a mixture of exponentials [33]. In this paper we adopt the hybrid Bernoulli/Gamma error metric proposed by Williams [32]. The distribution of the amount of precipitation, X , is modelled by

$$P(X > x) = \begin{cases} 1 & \text{if } x < 0 \\ \alpha \Gamma(\nu, \frac{x}{\theta}) & \text{if } x \geq 0 \end{cases} \quad (8)$$

where $0 \leq \alpha < 1$, $\nu > 0$, $\theta > 0$ and $\Gamma(\nu, z)$ is the (upper) incomplete Gamma function, $\Gamma(\nu, z) = \Gamma(\nu)^{-1} \int_z^{\infty} y^{\nu-1} e^{-y} dy$. The model is then trained to approximate the conditional probability of rainfall $\alpha(\mathbf{x}_i)$ and the scale, $\theta(\mathbf{x}_i)$, and shape, $\nu(\mathbf{x}_i)$, parameters of a Gamma distribution modelling the predictive distribution of the amount of precipitation. Logistic and exponential activation functions are used in output layer neurons to ensure that the distributional parameters satisfy their respective constraints.

III. EXPLOITING PREDICTIVE UNCERTAINTY

Environmental modellers are commonly interested in the impacts of extreme events, for example the impact of changes in future climate on local rainfall and subsequently on the flood hazard in susceptible catchments. General circulation models are considered to provide the best basis for estimating future climates that might result from anthropogenic modification of the atmospheric composition (i.e., the enhanced greenhouse effect). However, output from these models cannot be widely or directly applied in many impact studies because of their relatively coarse spatial resolution. The mismatch in scales between model resolution and the increasingly small scales required by impacts (e.g., agriculture and hydrology) analyses can be overcome by downscaling. Two major approaches to downscaling, statistical and dynamical (the latter using physically-based regional climate models), have been developed and tested in recent years, and shown to offer good potential for the construction of high-resolution scenarios of future climate change [10, 15, 28, 34]. Statistical downscaling methods seek to model the relationship between large scale atmospheric circulation, on say a European scale, and climatic variables, such as temperature and precipitation, on a regional or sub-regional scale, based on the historical record. Downscaling is an important area of research as it bridges the gap between predictions of future circulation generated by General Circulation Models (GCMs) and the effects of climate change on smaller scales, which are often of greater interest to end-users.

In order to estimate the impacts of changes in future climate on flood hazard, the predictions of a general circulation model are downscaled to provide predictions of future precipitation patterns, which in turn are processed by a

hydrological model to assess the effect of changes in rainfall patterns on water-levels in the river fed by the catchment being studied. In this example, we will consider a fictitious catchment¹ in which there is a flood hazard if the three-day total precipitation is in excess of 35 cm. Figure 1 shows a plot of the financial loss associated with flood events as a function of the three-day total precipitation; the loss is modelled as a constant component that is incurred whenever the river is unable to contain the run-off, and a component that reflects the additional damage resulting from increasingly severe flood events.

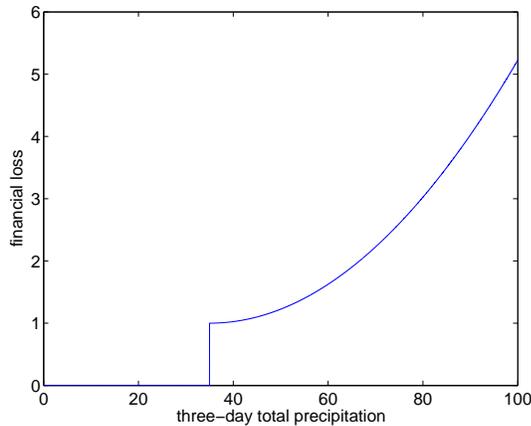


Fig. 1. Financial loss associated with flood events in a susceptible catchment as a function of the three-day total precipitation.

Figure 2 shows the three-day total precipitation time series for the study catchment area for the period 1979-1993. Note that many of the apparent dry spells are caused by missing data in the historical record rather than the absence of precipitation and are not included in the analysis. The measured loss for the observed time series is 49.02 units.

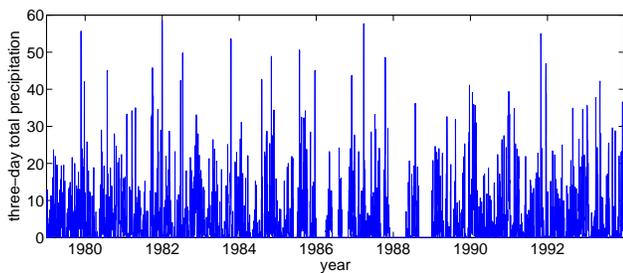


Fig. 2. Three-day total precipitation time series for a catchment area susceptible to flooding.

Figure 3 shows the predicted three-day total precipitation based on a conventional neural network downscaling model trained to estimate the conditional mean of the target distribution. The network was trained on two segments of the

precipitation time series spanning the periods 1961–1978 and 1994–2000. Note that the conditional mean systematically under-predicts the extreme rainfall events, as the predictive distribution is highly skewed. As a result, the predicted loss according to the simple neural network downscaling model is only 8.22 units, which severely under-estimates the true loss.

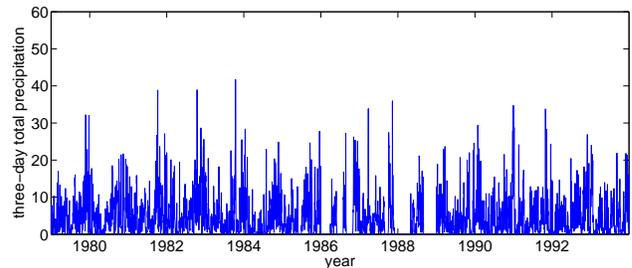


Fig. 3. Predicted three-day total precipitation time series for a catchment area susceptible to flooding, using a neural network providing the conditional mean of the target distribution.

A second neural network downscaling model was trained, this time using the hybrid Bernoulli/Gamma data misfit term (8). In this case, the model has three outputs, one supplies and estimate of the probability of rainfall and two that define a Gamma distribution modelling the plausibility of different amounts of rainfall. As this model provides a full probabilistic prediction, it is possible to generate synthetic precipitation time series, using the neural network as a conditional weather generator model. In order to infer the expected loss associated with the flood hazard, a Monte Carlo simulation is conducted using 100,000 synthetic precipitation time series generated by the network. Figure 4 shows a histogram of the measured losses from the Monte Carlo simulation, clearly the actual loss of 49.02 units is plausible, given the prediction distribution of loss. The expected loss, via Monte-Carlo integration, is 70.72 units, which is much closer to the recorded loss.

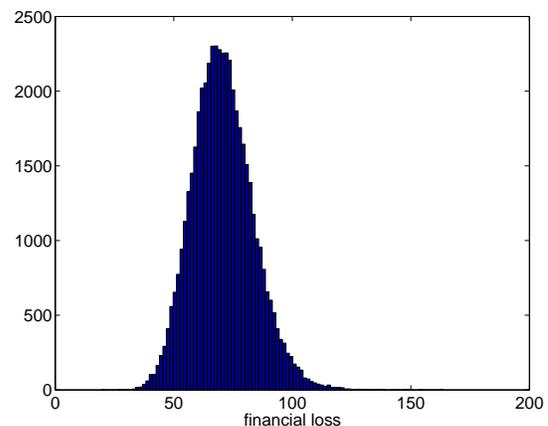


Fig. 4. Distribution of expected financial loss associated with flood events in a susceptible catchment.

¹The results are actually based on downscaled predictions for a real precipitation time series data from Newton Rigg, a rather wet station in the North West of the United Kingdom.

While this example is deliberately somewhat contrived, it does demonstrate that a probabilistic characterisation of the uncertainty of model predictions can be exploited in impact studies, especially where the principal focus lies on the implications of extreme events, which by their very nature are not modelled well by the conditional mean. The integration over sources of uncertainty also provides the results in a format that is well suited to the needs of end-users, such as government institutions or the insurance industry. Clearly the distribution of plausible losses is exactly the information required by such users for well-informed policy-making and forward planning.

IV. THE PREDICTIVE UNCERTAINTY IN ENVIRONMENTAL MODELLING CHALLENGE

The WCCI-2006 predictive uncertainty in environmental modelling challenge consisted of one SYNTHETIC benchmark dataset and three real-world environmental datasets PRECIP, SO2 and TEMP. The format of the competition was based closely on the regression problems of the earlier Pascal predictive uncertainty challenge. The negative log-likelihood of the test data was used as the performance criterion for the final ranking of submissions, as it is the natural measure of the fit of a distribution to a set of data. Two standard methods were available for describing the predictive distribution for each pattern, the mean and variance of a Gaussian predictive distribution, or a set of quantiles, allowing the definition of arbitrary predictive distribution. An unusual feature of the competition is that the competitors had the option of suggesting alternate forms for specifying the predictive distribution (as the likelihood can be described in any number of parametric forms). A mixture Gaussian option was added at a late stage in the competition in response to a request from one of the competitors. The target data for all three of the real-world environmental benchmark datasets are (finely) quantised, for example precipitation data is only measured to the nearest 0.1 mm. In principle it would therefore be possible to make the negative log-likelihood arbitrarily low by specifying the predictive distribution (via quantiles) as a set of delta functions centred on the quantised values. This technique was employed by some entries to the original Pascal predictive uncertainty challenge. In order to prevent this, the minimum allowable width of the quantiles (and similarly the variances of the individual components of a mixture Gaussian predictive distribution) were limited to match the quantisation interval used.

A. Reference Submissions

Three baseline models were submitted for each dataset, which gave a fixed predictive distribution for all patterns: Baseline #1 - fixed Gaussian predictive distribution specified via the unconditional mean and variance of the target data, Baseline #2 - fixed Gaussian distribution specified as a set of quantiles and Baseline #3 - fixed predictive distribution specified by quantiles representing the empirical distribution of the target data. A fourth baseline model was

created for the SO2 benchmark, giving a fixed predictive distribution for all patterns based on a Gaussian mixture model of five components, fitted using the standard Expectation Maximisation (EM) algorithm, as implemented by the NET-LAB package. In addition to these baseline models, neural network models were also submitted for each benchmark, the training procedure used is described in Section II. A heteroscedastic Gaussian data mis-fit term (7) was used for the SYNTHETIC and TEMP benchmarks, a heteroscedastic log-normal term for the SO2 benchmark and the hybrid Bernoulli/Gamma term (8) term for the PRECIP benchmark. In order to avoid training difficulties due to local minima of the cost function, 20 models were trained in each case, with randomly initialised weights, and the model giving the lowest value for the regularised loss retained. These models provide an indication of the “minimum” and “competitive” levels of performance for each benchmark.

B. The SYNTHETIC Benchmark

A synthetic heteroscedastic regression problem, taken from Williams [31], was included in the challenge, principally to provide a relatively small dataset that could be easily visualised for the purposes of initial model development. However, as the true conditional mean and variance functions are known, it is straight-forward to assess the quality of the model. The univariate input patterns, x , are drawn from a uniform distribution on the interval $(0, \pi)$, the corresponding targets, y , are drawn from a univariate Normal distribution with mean and variance that vary smoothly with x :

$$x_i \sim \mathcal{U}(0, \pi),$$

$$y_i \sim \mathcal{N}\left(\sin\left[\frac{5x}{2}\right] \sin\left[\frac{3x}{2}\right], \frac{1}{100} + \frac{1}{4} \left[1 - \sin\left[\frac{5x}{2}\right]\right]^2\right).$$

Figure 5 shows a plot of the synthetic benchmark dataset, along with indications of the true conditional mean and standard deviation. The heteroscedastic (input-dependent variance) nature of the data is clearly evident.

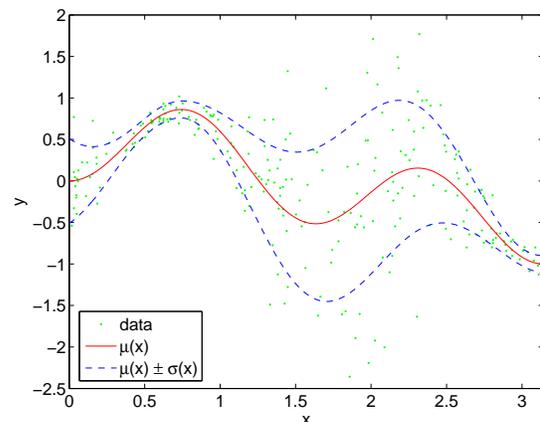


Fig. 5. Plot of the training data for the SYNTHETIC benchmark dataset, along with an indication of the true conditional mean, $\mu(x)$ and conditional standard deviation, $\sigma(x)$.

TABLE I

TRAINING AND TEST SET NEGATIVE LOG-LIKELIHOOD STATISTICS FOR THE SYNTHETIC BENCHMARK.

Name	Method	Train NLPD	Test NLPD
Reference	ground truth	0.3333	0.3489
Harva	varmlp (MoG)	0.3251	0.3858
Cawley	MLP	0.3083	0.4046
Kurogi <i>et al.</i>	CAN2 ensemble + CV	0.2236	0.4304
Boardman	Support Vector Regression	0.4150	0.4745
Nikulin	CM+GbO	0.3590	0.4805
Bagnall	YJ	1.0081	1.0313
Reference	Baseline #1	1.1064	1.1357
Reference	Baseline #2	1.1104	1.1374
Reference	Baseline #3	0.7923	1.2324

Table I shows the negative log-likelihood of the training and test sets of the SYNTHETIC benchmark for selected entries. It can be seen that many of the entries were able to make clear improvements in modelling the predictive distribution over the baseline models, with the best models approaching the performance of the optimal “ground truth” model used to generate the data. However, the SYNTHETIC benchmark is relatively straight-forward, the only unusual feature being the heteroscedasticity of the noise process.

C. The PRECIP Benchmark

The PRECIP benchmark models a realistic statistical downscaling exercise, the aim of which is to predict the (scaled) precipitation for Newton Rigg, a relatively wet station in the North-West of the United Kingdom, using inputs representing large scale circulation features (see [6, 14] for further details). Figure 6 shows a histogram of the target data for the training set of the PRECIP benchmark, highlighting a number of unusual features of this dataset. Firstly, the data is non-negative (it would make little sense to talk of negative rainfall). Secondly, there is a large probability mass centred on zero, representing the proportion of days where no rainfall occurs. Rainfall presents an example of a *mixed* distribution, and is often modelled as separate occurrence and amount processes, where the probability of rainfall is given by, e.g. a logistic regression model, and the amount of rainfall given by a e.g. linear regression model fitted to the training data representing days where rainfall was actually observed. The hybrid Bernoulli/Gamma mis-fit term (8) simply combines the occurrence and amount processes as a single model within the framework of maximum likelihood. The extra probability mass at the origin is easily accommodated by the quantile representation of the predictive distribution.

Table II shows the test set MSE and NLPD statistics for selected models. Note that while many model were able to make useful reductions in the mean-squared error, only one of the submissions was able to improve on Baseline #1. This suggests that current modelling techniques are likely to be inadequate for use in statistical impact studies of the

TABLE II

TEST SET MEAN-SQUARED ERROR (MSE) AND NEGATIVE LOG-LIKELIHOOD (NLPD) STATISTICS FOR THE PRECIP BENCHMARK.

Name	Method	MSE	NLPD
Cawley	MLP	0.6305	-0.5095
Harva	varmlp	5.4493	-0.2792
Reference	Baseline #1	1.0002	-0.1772
Takeuchi	Kernel QR	0.6109	0.7469
Bagnall	YJ	2.1072	1.1139
Nikulin	CM+GbO	0.6539	1.2724
Boardman	Support Vector Regression	0.6441	1.6055
Reference	Baseline #2	1.0001	2.0346
Reference	Baseline #3	1.0001	2.0496
Kurogi <i>et al.</i>	CAN2 ensemble + CV + hetero + quantile	0.6465	3.0982

nature described in Section III.

D. The SO2 Benchmark

The SO2 benchmark represents an atmospheric pollution forecasting problem, where the aim is to predict 24 hours in advance the SO₂ concentration in urban Belfast, based on meteorological conditions and current SO₂ levels (see [23] for further details). Table III shows the test set mean-squared error and negative log likelihood for selected models over the SO2 benchmark. Clearly this is the noisiest of the benchmark datasets, and while some reduction in the mean-squared-error is possible, it is difficult to produce a model that improves on the baseline models in terms of the quality of the predictive distribution.

E. The TEMP Benchmark

The TEMP benchmark problem is perhaps the most easily modelled of the real-world benchmark problems, and again represents a downscaling problem, where the aim in this case is to model the daily maximum temperature at the Writtle station in the South-East of the United Kingdom base on

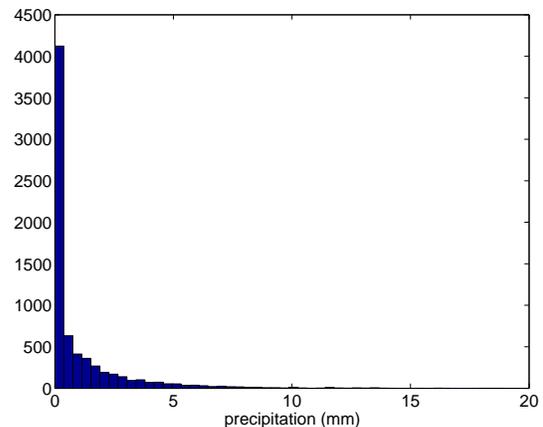


Fig. 6. Histogram of the target data for the training set of the PRECIP benchmark dataset.

TABLE III

TEST SET MEAN-SQUARED ERROR (MSE) AND NEGATIVE LOG-LIKELIHOOD (NLPD) STATISTICS FOR THE SO₂ BENCHMARK.

Name	Method	MSE	NLPD
Cawley	MLP	0.7985	4.2550
Harva	varmlp	0.8333	4.3702
Reference	Baseline #4	1.0000	4.4964
Reference	Baseline #1	1.0001	4.4968
Nikulin	CM+GbO	0.8576	4.6162
Bagnall	YJ	1.7598	4.7578
Boardman	Support Vector Regression	0.8396	5.0897
Reference	Baseline #3	1.0000	5.1655
Reference	Baseline #2	1.0000	5.2181
Takeuchi	Kernel QR	0.6884	6.0425
Kurogi <i>et al.</i>	CAN2 ensemble + CV + hetero + quantile	0.7807	11.0063

TABLE IV

TEST SET MEAN-SQUARED ERROR (MSE) AND NEGATIVE LOG-LIKELIHOOD (NLPD) STATISTICS FOR THE TEMP BENCHMARK.

Name	Method	MSE	NLPD
Snelson	Sparse pseudo-input Gaussian process (SPGP)	0.0661	0.0348
Cawley	MLP	0.0693	0.0530
Kurogi <i>et al.</i>	CAN2 ensemble + CV + hetero + quantile + outlier	0.0681	0.0591
Boardman	Support Vector Regression	0.0709	0.0760
Nikulin	CM+GbO	0.0729	0.1076
Bagnall	Linear Regression	0.077432	0.136235
Harva	varmlp	0.0925	0.2015
Whittlely	QuantLin	24.9839	0.6251
Reference	Baseline #1	1.0000	1.3004
Reference	Baseline #2	1.0000	1.4151
Reference	Baseline #3	1.0000	1.4177
Takeuchi	Kernel QR	0.0965	24.7922

similar large scale circulation features as those used for the PRECIP benchmark. In this case a heteroscedastic Gaussian noise process is a reasonable assumption. Table IV shows the test set MSE and NLPD statistics for selected models, in almost all cases the models significantly improve on the baseline models in terms of the NLPD.

F. Final Competition Standings

The final standings in the competition, decided by mean NLPD score over the three environmental datasets, are shown in Table V. The overall winner is Markus Harva.

V. AREAS FOR FURTHER RESEARCH

A. Inherent Bias in the Conditional Variance

It is well known that estimates of the conditional variance are likely to be significantly biased. If the model of the conditional mean over-fits the data, this reduces the apparent local noise density, and so error bars based on the conditional variance will be unrealistically narrow. This problem has

TABLE V

FINAL STANDINGS IN THE COMPETITION - THE OVERALL WINNER, DECIDED BY MEAN NLPD SCORE, IS MARKUS HARVA.

Name	PRECIP	SO ₂	TEMP	Mean
Cawley	-0.5095	4.2550	0.0530	1.2661
Harva	-0.2792	4.3702	0.2015	1.4308
Nikulin	1.2724	4.6162	0.1076	1.9987
Bagnall	1.1139	4.7578	0.1362	2.0026
Boardman	1.6055	5.0897	0.0760	2.2571
Kurogi <i>et al.</i>	3.0982	11.0063	0.0591	4.7212
Takeuchi	0.7469	6.0425	24.7922	10.5272
Whittlely	∞	∞	0.6251	∞
Snelson	∞	∞	0.0348	∞

previously been addressed via Bayesian approaches [4, 11], and by the use of leave-one-out cross-validation [7]. However these approaches are currently only suitable for relatively small scale applications, with only a few thousands of training patterns. Further research is needed to develop large scale algorithms suitable for environmental applications, where much larger amounts of data are typically available.

B. Incorporating the Uncertainty in the Model Parameters

In this paper we have reviewed the use of maximum-likelihood based loss functions for neural networks, which allow us to incorporate prior knowledge regarding the uncertainty in model predictions due to the inherent noise process contaminating the data. Another important source of uncertainty lies in the uncertainty due to the estimation of the model parameters from a finite sample of data. It seems likely that a better model of the predictive distribution might be obtained by including this effects of the uncertainty in the model parameters, e.g. via the Laplace approximation [17, 18] or via Markov-Chain Monte Carlo methods [20].

C. The Form of the Predictive Distribution

While expert knowledge is sometimes available regarding the form of the noise process contaminating the data, it would be useful also to have a data-driven approach, where the form of the noise process is also inferred from the training data. The mixture density network [2], where the output of the model specify the components of a Gaussian mixture model of the predictive distribution, represents the most basic approach. The warped Gaussian Process, [25], in which the observation space is transformed so as to be well modelled as a Gaussian process, represents a more recent approach.

VI. CONCLUSIONS

In this paper we have demonstrated that a model of the predictive distribution can be exploited in studies of the impacts of changes in future climate, via a somewhat contrived, but nevertheless illustrative example. An online competition has been organised in an attempt to promote research on methods for estimating the uncertainty inherent in statistical predictions. The results demonstrate that this

is a difficult topic, where standard approaches do not yield uniformly good results. We hope that the competition has gone some way to highlight an area where further research is likely to produce practical benefits in the analysis of environmental data.

ACKNOWLEDGEMENTS

Many thanks are due to Joaquin Quinero Candela, Carl Rasmussen and Yoshua Bengio for organising the original (and very successful) Evaluating Predictive Uncertainty Challenge, and to Vladimir Cherkassky, Julio Valdes, Vladimir Krasnopolsky and Dimitri Solomatine for organising the special session on Applications of Learning and Data-Driven Methods to Earth Sciences and Climate Modeling at IJCNN-2005. These two events provided the impetus (and the format) for this competition. Thanks to the anonymous reviewers for their helpful comments. Thanks are also due to Radford Neal, Iain Murray and Markus Harva for their helpful feedback on the format of the competition and to Peter Williams for provoking our initial interest in modelling predictive distributions. Thank you to the anonymous reviewers for their helpful and constructive comments on a previous draft of this paper. The competition would not have been possible without Nicola Talbot, who implemented the competition website. Lastly, many thanks to all those who took part!

REFERENCES

- [1] B. Bhattacharya and D. P. Solomatine. Neural networks and M5 model trees in modelling water level-discharge relationship. *Neurocomputing*, 63:381–396, January 2005.
- [2] C. M. Bishop. Mixture density networks. Technical Report NCRG/94/004, Neural Computation Research Group, Aston University, 1994.
- [3] C. M. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, 1995.
- [4] C. M. Bishop and C. S. Qazaz. Bayesian inference of noise levels in regression. In C. von der Malsburg, W. von Seelen, J. C. Vorbrüggen, and B. Sendhoff, editors, *Proceedings of the International Conference on Artificial Neural Networks (ICANN-96)*, volume 1112 of *Lecture Notes in Computer Science*, pages 59–64, Bochum, Germany, July 16–19 1996. Springer.
- [5] W. L. Buntine and A. S. Weigend. Bayesian back-propagation. *Complex Systems*, 5:603–643, 1991.
- [6] G. C. Cawley, S. R. Dorling, P. D. Jones, and C. Goodess. Statistical downscaling with artificial neural networks. In *Proceedings of the European Symposium on Artificial Neural Networks (ESANN-2003)*, pages 167–172, Bruges, Belgium, April 23–25 2003.
- [7] G. C. Cawley, N. L. C. Talbot, R. J. Foxall, S. R. Dorling, and Mandic D. P. Heteroscedastic kernel ridge regression. *Neurocomputing*, 57:105–124, March 2004.
- [8] C. W. Dawson, L. M. See, R. J. Abraham, R. L. Wilby, A. Y. Shamseldin, F. Anctil, A. N. Belbachir, G. Bowden, G. Dandy, N. Lauzon, and H. Maier. A comparative study of artificial neural network techniques for river stage forecasting. In *Proceedings of the IEEE International Joint Conference on Neural Networks (IJCNN '05)*, volume 4, pages 2666–2670, 31 July - 4 August 2005.
- [9] S. Geman, E. Bienenstock, and R. Doursat. Neural networks and the bias/variance dilemma. *Neural Computation*, 4(1):1–58, 1992.
- [10] F. Giorgi and L. O. Mearns. Introduction to special section: Regional climate modeling revisited. *Journal of Geophysical Research*, 104:6335–6352, 1999.
- [11] P. Goldberg, C. Williams, and C. Bishop. Regression with input-dependent noise: A Gaussian process treatment. In M. Kearns, M. Jordan, and S. Solla, editors, *Advances in Neural Information Processing Systems*, volume 10, pages 493–499. MIT Press, Cambridge, MA, 1998.
- [12] I. S. Gradshteyn and I. M. Ryzhik. *Table of Integrals, Series and Products*. Academic Press, fifth edition, 1994.
- [13] C. Harpham and R. L. Wilby. Multi-site downscaling of heavy daily precipitation occurrence and amounts. *Journal of Hydrology*, 312(1–4):235–255, October 2005.
- [14] M. R. Haylock, G. C. Cawley, C. Harpham, R. L. Wilby, and C. Goodess. Downscaling heavy precipitation over the UK: A comparison of dynamic and statistical methods and their future scenarios. *International Journal of Climatology* (in press), 2006.
- [15] B. C. Hewitson and R. G. Crane. Climate downscaling: Techniques and application. *Climate Research*, 7:85–95, 1996.
- [16] H. S. Jeffreys. *Theory of Probability*. Oxford University Press, 1939.
- [17] D. J. C. Mackay. Bayesian interpolation. *Neural Computation*, 4(3):415–447, 1992.
- [18] D. J. C. Mackay. A practical Bayesian framework for backprop networks. *Neural Computation*, 4:448–472, 1992.
- [19] D. J. C. MacKay and Z. Gharahmani. Comments on ‘maximum likelihood estimation of intrinsic dimension’ by Levina and Bickel. <http://www.inference.phy.cam.ac.uk/mackay/dimension>, 2005.
- [20] R. M. Neal. *Bayesian Learning for Neural Networks*. Springer, New York, 1996.
- [21] D. A. Nix and A. S. Weigend. Estimating the mean and variance of the target probability distribution. In *Proc., Int. Conf. on Neural Networks*, volume 1, pages 55–60, 1994.
- [22] D. A. Nix and A. S. Weigend. Learning local error bars for nonlinear regression. In *Advances in Neural Information Processing Systems*, volume 7, pages 489–496. MIT Press, 1995.
- [23] G. Nunnari, S. R. Dorling, U. Schlink, G. Cawley, R. Foxall, and T. Chatterton. Modelling SO₂ concentration at a point with statistical approaches. *Environmental Modelling and Software*, 19(10):887–905, October 2004.
- [24] U. Schlink, S. Dorling, E. Pelikan, G. Nunnari, G. Cawley, H. Junninen, A. Greig, R. Foxall, K. Eben, T. Chatterton, J. Vondracek, M. Richter, M. Dostal, L. Bertuccio, M. Kolehmainen, and M. Doyle. A rigorous inter-comparison of ground-level ozone predictions. *Atmospheric Environment*, 37(23):3237–3253, July 2003.
- [25] E. Snelson, Rasmussen C. E., and Z. Ghahramani. Warped gaussian processes. In S. Thrun, L. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing Systems*, volume 16, pages 337–344. MIT Press, Cambridge, MA, 2004.
- [26] R. D. Stern and R. Coe. A model fitting analysis of daily rainfall data (with discussion). *Journal of the Royal Statistical Society A*, 147(1):1–34, 1984.
- [27] A. N. Tikhonov and V. Y. Arsenin. *Solutions of ill-posed problems*. John Wiley, New York, 1977.
- [28] R. L. Wilby, T. M. L. Wigley, D. Conway, P. D. Jones, B. C. Hewitson, J. Main, and D. S. Wilks. Statistical downscaling of general circulation model output: A comparison of methods. *Water Resources Research*, 34:2995–3008, 1998.
- [29] P. M. Williams. A Marquardt algorithm for choosing the step-size in backpropagation learning with conjugate gradients. Cognitive Science Research Paper CSRP-229, University of Sussex, Brighton, U.K., February 1991.
- [30] P. M. Williams. Bayesian regularisation and pruning using a Laplace prior. *Neural Computation*, 7(1):117–143, 1995.
- [31] P. M. Williams. Using neural networks to model conditional multivariate densities. *Neural Computation*, 8:843–854, 1996.
- [32] P. M. Williams. Modelling seasonality and trends in daily rainfall data. In M. I. Jordan, M. J. Kearns, and S. A. Solla, editors, *Advances in Neural Information Processing Systems - Proceedings of the 1997 Conference*, volume 10, pages 985–991. MIT Press, 1998.
- [33] D. A. Woolhiser and G. G. S. Pegram. Maximum likelihood estimation of Fourier coefficients to describe seasonal variation of parameters in stochastic daily precipitation models. *Journal of Applied Meteorology*, 18:34–42, 1979.
- [34] E. Zorita and H. von Storch. The analog method as a simple statistical downscaling technique: Comparison with more complicated methods. *Journal of Climate*, 12:2474–2489, 1999.