Machine learning techniques in predicting uncertainty of environmental models

#### **Dimitri Solomatine**

Professor of Hydroinformatics, IHE Delft Institute for Water Education Delft, The Netherlands



## Outline

- Introduction: what are analysisng?
- Machine learning methods to (a) analyse and (b) predict the model uncertainty
- Suggested approach: "escalation" of uncertainty
- Examples

## **Example for a quick start: deterministic forecasts and 90% uncertainty bounds**



## Sources of model uncertainty: perceptual, structure, parameters, data

 $y = M(x, p) + \varepsilon_s + \varepsilon_{\theta} + \varepsilon_x + \varepsilon_v$ 



D.P. Solomatine. Escalation of uncertainty.

## Traditional steps in uncertainty analysis of a calibrated model

- Identification of sources of uncertainty (input, parameter, model structure)
- Quantification of uncertainty (e.g. as distribution)
- Studying *propagation of uncertainty* through the model (e.g. by Monte Carlo simulation)
- Quantification of uncertainty in the model outputs (i.e. identification of output distribution (pdf) or its characteristics mean, st.dev., quantiles)
- If possible, reduction of uncertainty (e.g. model improvement, more accurate measurements, etc.)
- Application of the uncertain information in decision making process

## Data uncertainty (input, parameters): propagation of uncertainty through the model

- y^ = M (x, p)
- x = input, p = parameters
  - Uncertainty in X and p *propagates* to output y
  - pdf of parameters  $\rightarrow$  pdf of output
  - pdf of inputs  $pdf_x \rightarrow pdf$  of output  $pdf_x \rightarrow pdf_y$
- $pdf_p \rightarrow pdf_y$  $pdf_x \rightarrow pdf_y$

### **Monte Carlo Simulation**

### Mote Carlo casino: roulette wheel







 It is a random number generator – uses uniform distribution with the range of [0, 36]





## Monte Carlo simulation in analysing parametric uncertainty



#### Sampling parameters and multiple model runs





#### Sampling rainfall and multiple model runs



## **Representing uncertainty of model output by the confidence bounds**



Instead of fitting a theoretical distribution, we can use mean, standard deviation, quantiles.

E.g., 5% and 95% form the **90% confidence bounds** 

Propagation of parameters/data uncertainty by Monte Carlo simulation is a typical practical approach.

But is it the only one?

## QUESTION 1. On assumptions

- We are assuming some known distributions of parameters or inputs. How safe is this?
- Could we take a safer route and assume less?
- Let's make a step backwards and pose the

**QUESTION 1:** 

what is the uncertainty of the calibrated model itself?

# **Residual uncertainty: uncertainty of a calibrated ("optimal") model**



- Uncertainty of an *optimal model* M  $(x, \theta)$ 
  - Model M is calibrated on measured data y
  - We say the model M uncertainty is manifested in the residual model error  $\varepsilon = y^{-}y$
  - This error incorporates all uncertainties due to: observational errors, inaccurately estimated parameters, inadequate model structure

D.P. Solomatine. Escalation of uncertainty.

## ESCALATION ("build up") of model uncertainty [message 1]

- 1. Study the (residual) uncertainty of an *optimal model* M (p\*)
- 2. Add and study (typically, by MC simulation)
  - A) uncertainty of M (p\*) due to DATA uncertainty
  - B) uncertainty of M (p) due to PARAMETERS uncertainty
- 3. Add and study uncertainty of M (p) due to STRUCTURAL uncertainty
  - 4. Study uncertainty of a *model class* M (p), given the probabilistic properties of parameters and data

## QUESTION 2. On what is analysed



- In UA we always use the past data, so Estimates of uncertainty are about the PAST.
- QUESTION 2:

*how can we assess the model uncertainty for new inputs, i.e. for the future?* 

- and this question we pose for all sources of uncertainty (and not only residual)

D.P. Solomatine. Escalation of uncertainty.

Models of Residual Uncertainty : Using Methods of Computational Intelligence

#### CI in building models of natural processes why not build a model of uncertainty?



- CI provides methods to build *Data-driven* models
- Ideally, such models are "ultimate models" since they are not polluted by theories

D.P. Solomatine. Escalation of uncertainty.

## Example of a data-driven (statistical, CI) model

- observed data characterises the input-output relationship
  X → Y
- model parameters are found by optimization
- the model then predicts output for the new input without actual knowledge of *what* drives Y



Which model is "better": green, red or blue?

## **CI models: are they indeed intelligent?**

#### Artificial neural network



$$Y = g^{out} (b_{0k} + \sum_{j} b_{jk} g^{hid} (a_{0j} + \sum_{i} a_{ij} x_{i}^{(t)}))^{2} \quad where \quad g(u) = \frac{1}{1 + e^{-\alpha u}}$$

D.P. Solomatine. Escalation of uncertainty.

## Data-driven model as an *error corrector* for a process (physically-based) model



# Data-driven model to predict the residual error distribution



Train data-driven model (e.g. Neural Network) to forecast residual error *pdf* (*i.e. the model* **M** *output uncertainty*)

# Some of the models of residual uncertainty

- QR (1978) (quantile regression): autoregressive linear model of model residuals predicts the distribution quantiles [Koenker & Basset]
- DUMBRAE (2012) (Dynamic Uncertainty Model By Regression on Absolute Error) [Pianosi & Raso]: autoregressive model of model residuals (it corrects the model residual first and then carries out the uncertainty prediction by an autoregressive model)

UNNEC (2006, 2009) (UNcertainty Estimation based on local Errors and Clustering) [Shrestha & Solomatine]: it takes into account all variables influencing such uncertainty and uses machine learning (non-linear) methods (neural networks, model trees, instance-based learning etc.)

D.P. Solomatine. Escalation of uncertainty.

## **UNEEC** method

UNcertainty Estimation based on local Errors and Clustering

 machine learning model of the *past residual errors of the* optimal process model is built

D.P. Solomatine, D.L. Shrestha (2009). A novel method to estimate model uncertainty using machine learning techniques. *Water Resources Res.* 45, W00B11.

## **UNEEC:** assumptions, constraints

#### Assumptions

- Model error is an indicator of the model uncertainty
- Model error depends on the current condition of a natural system and can be predicted
- Model errors are similar for similar conditions
- Constraints
  - Model structure and parameters are fixed
  - Need to re-train the error model with the changes in the catchment characteristics (e.g. land use change)
  - Data hungry, more data are needed for reliable results

## Idea 1: assess CDF from all historical data about errors



## Idea 2: local modelling of errors (link CDF to "characteristis variables")



## Idea 3: Use fuzzy clustering of examples to generate training data sets



#### **Using instance-based learning**



#### **UNEEC details. Step 1: clustering**

Clustering (finding groups of data in the space characterising hydro-meteo condition): K-means clustering, fuzzy C-means clustering

**Obj. function**  $\min(U,V) \left\{ J_m(U,V) = \sum_{j=1}^{c} \sum_{i=1}^{N} \mu_{i,j}^m D_{i,j}^2 \right\}$ 

 $D_{i,j}^2 = \|x_i - v_j\|_{A}^2$ 

Distance

Degree of  $m \ge 1$ Fuzzification

Constraint

$$\sum_{j=1}^{c} \mu_{i,j} = 1, \forall i$$



#### **UNEEC details. Step 2: Determining Prediction Interval (PI) for each cluster**



D.P. Solomatine. Escalation of uncertainty.

#### UNEEC details. Step 3, 4, 5: Building and using the model



## **UNEEC methodology**



## Extensions (simplifications) of UNEEC: without clustering and using instance-based learning



Based on Master study of Omar Wani (2015)

- SKIBLUE (Streamflow-Centric K nearest neighbour Instance-Based Learning and Uncertainty Estimation)
- O. Wani, J. Beckers, A.H. Weerts, D.P. Solomatine. Nonparametric Predictive Uncertainty Estimation Using Instance Based Learning with Applications to Hydrologic Forecasting. HESS-D, 2016.
- Based on Master study of Ms. Jingyi Chen (2015)
  - UNEEC-IBL
  - Jingyi Chen (2015). Uncertainty Prediction in Hydrological Modelling: Case of Dapoling-Wangjiaba Catchment in Huai River Basin. UNESCO-IHE Master thesis D.P. Solomatine. Escalation of uncertainty.

— Escalating uncertainty —

Assuming now uncertainty in parameters and or data...

Running Monte Carlo simulations...

But how to estimate output uncertainty for the new model runs?

Models of Parametric Uncertainty : ... and again using Methods of Computational Intelligence

## **MLUE method**

#### Machine Learning in Uncertainty Estimation

 machine learning model of the *process model's Monte Carlo simulation results* is built

D. L. Shrestha, N. Kayastha, and D. P. Solomatine (2009). A novel approach to parameter uncertainty analysis of hydrological models using neural networks. *HESS*, 13, 1235–1248.

#### **Monte Carlo simulation of parametric uncertainty**



# Issues with MC for new model runs in real-time

#### Issues with re-running MC for new inputs:

- 1) convergence of the Monte Carlo simulation is very slow (O(N^-0.5)) so larger number of runs needed to establish a reliable estimate of uncertainties
- 2) number of simulation increases exponentially with the dimension of the parameter vector ((O(n^d)) to cover the entire parameter domain

#### Idea:

 encapsulate the results of MC simulation in a machine learning model

## **MLUE Methodology (1)**

- Consider the sources of the uncertainty analysis to be conducted within the framework of Monte Carlo simulation
- Execute the MC simulations to generate the data y<sub>i</sub>(t) = M (X(t), p<sub>i</sub>)
- Estimate the uncertainty measures of the MC realizations, e.g., mean, variance, prediction intervals, quantiles
  - to start with, estimate two quantiles (say, 5% and 95%), forming the prediction interval PI

## **MLUE Methodology (2)**

- Analyze the dependency of the uncertainty measures (quantiles) on the input and state variables of the hydrological model
  - we used Correlation and Average mutual information analysis
- Select the input variables for machine learning model based on the dependency analysis
- Train the machine learning model U to predict the uncertainty measures of MC realizations  $PI = U(\mathbf{X})$
- Validate machine learning model U by estimating the uncertainty measures with the "new" input data
- Use model U

## Validation

- Measuring predictive capability of uncertainty model U (measures the accuracy of uncertainty models in approximating the quantiles of the model outputs generated by MC simulations)
  - Coefficient of correlation (r) and root mean squared error (RMSE)
- Measuring the statistics of the uncertainty estimation (i.e. goodness of the model U as uncertainty estimator)
  - Prediction interval coverage probability (PICP) and mean prediction interval (MPI) (Shrestha & Solomatine 2006, 2008)

$$PICP = \frac{1}{n} \sum_{t=1}^{n} C$$
  
with  $C = \begin{cases} 1, \ PL_t^L \le y_t \le PL_t^U \\ 0, \ \text{otherwise} \end{cases}$   $MPI = \frac{1}{n} \sum_{t=1}^{n} (PL_t^U - PL_t^L)$ 

 Visualizing such as scatter and time plot of the prediction intervals obtained from the MC simulation and their predicted values



### **Applications**

UNEEC and MLUE were tested and compared to other methods on 5 various cases: *Brue, Bagmati, Sieve, Severn, Dapoling-Wanjiaba* 

### **Study area: Brue catchment, UK**



## Study area: Brue catchment, UK



## **Conceptual Hydrological model HBV**



## **Data Analysis**

- Analysis of dependency btw various combinations of the input variables and the output
  - Correlation
  - Average mutual information (AMI) between *REt* and PIs, ( optimal lag time is around 7-9 hours).
  - Additional analysis of the correlation and AMI between the PIs and observed discharge *Qt* are carried out. (i.e. with the lag of 0, 1, 2) have very high correlation with the PIs.

## **Experimental setup**

#### MC simulation

- 9 Parameters of HBV model are sampled uniformly from the feasible ranges
- Nash-Sutcliffe coefficient of efficiency (CE) is used as error measure
- Convergence stabilized after 10,000 (75,000 runs made)
- Only 25,000 "good" models considered (rejection threshold is set to 0) to compute prediction quantiles

## **Experimental setup**

#### Machine learning model U

- $\bullet PI = U(RE_{t-5a'} Q_{t-1'} \Delta Q_{t-1})$ 
  - PI lower or upper prediction intervals,
  - $RE_{t-5a}$  average of  $RE_{t-5r}$   $RE_{t-6r}$   $RE_{t-7r}$   $RE_{t-8r}$  and  $RE_{t-9}$

• 
$$\Delta Q_{t-1} - Q_{t-1} - Q_{t-2}$$

 Input variables were selected based on the analysis of their relatedness to output error (average mutual information)

$$AMI = \sum_{i,j} P_{XY}(x_i, y_j) \log_2 \left[ \frac{P_{XY}(x_i, y_j)}{P_X(x_i) P_Y(y_j)} \right]$$

- Methods:
  - M5 model trees,
  - Iocally weighted regression
  - MLP neural networks

### **Results**

## **UNEEC: Clustering result example**



D.P. Solomatine. Escalation of uncertainty.

## **UNEEC: Performance (MLP ANN)**



## **UNEEC: Estimation of prediction intervals**



## **MLUE: Estimation of prediction intervals**



## **MLUE: Performances**

#### Predictive capability

	Corr C		RMSE	
	PIL	<b>ΡΙ</b> υ	PI <sup>∟</sup>	PI∪
МТ	0.841	0.792	0.614	1.641
LWR	0.822	0.798	0.643	1.604
ANN	0.847	0.806	0.584	1.568

#### Goodness of uncertainty measures

	MCS	MT	LWR	ANN
PICP %	77.24	66.97	75.16	65.54
MPI m³/s	2.09	2.03	1.93	1.96

MCS = Monte Carlo MT = M5 Model tree LWR = local weighted regression ANN =MLP neural network

## **Extensions**

- Estimation of several quantiles 5%, 10%:10%:90%, 95%
  - i.e. estimating cdf of MC realizations by machine learning models



D.P. Solomatine. Escalation of uncertainty.

## Use of Machine learning methods: conclusions

- Machine learning methods are able to replicate:
  - Past performance of a process model
  - Results of Monte-Carlo simulations
- The methods are computationally efficient and can be used in real time application
- They are to various kinds of models
- The results demonstrate that the interpretable uncertainty estimates are generated
- Future work:
  - Other ML methods are to be tested
  - The methods can be applied in the context of other sources of uncertainty - input, structure, or combined

## References

#### UNEEC and extensions:

- D.L. Shrestha, D.P. Solomatine. Machine learning approaches for estimation of prediction interval for the model output. *Neural Networks*, 2006, 19(2), 225-235.
- D.P. Solomatine, D.L. Shrestha. A novel method to estimate model uncertainty using machine learning techniques. *Water Resour Res.* 45, W00B11, 2009.
- N. Dogulu, P. López López, D. P. Solomatine, A. H. Weerts, and D. L. Shrestha. Estimation of predictive hydrologic uncertainty using the quantile regression and UNEEC methods and their comparison on contrasting catchments, *Hydrol. Earth Syst. Sci.*, 19, 3181-3201, 2015.
- O. Wani, J. Beckers, A.H. Weerts, D.P. Solomatine. Non-parametric Predictive Uncertainty Estimation Using Instance Based Learning with Applications to Hydrologic Forecasting. HESS-D, 2016

#### MLUE:

- D. L. Shrestha, N. Kayastha and D. P. Solomatine. A novel approach to parameter uncertainty analysis of hydrological models using neural networks. *Hydrol. Earth Syst. Sci.*, 13, 1235-1248, 2009.
- Shrestha, D.L., Kayastha, N., Solomatine, D., Price, R. Encapsulation of parametric uncertainty statistics by various predictive machine learning models: MLUE method. *J Hydroinformatics*, 16 (1), 95-113, 2014.

## Conclusions

- Uncertainty analysis should always contain explicit answers to two questions:
  - 1) what type of uncertainty is to be analysed: residual (which do not need MC), or parametric/data (which need MC)
  - 2) what is required: just analysis of the past, or also a model predicting the future uncertainty
- It is advisable:
  - to go explicitly through all stages of uncertainty escalation, starting from residual uncertainty
  - to try to build the **predictive models of uncertainty** at all stages
  - complement the deterministic models M with a *family of uncertainty* models U

## What to know more?

- We teach Master courses:
  - Hydroinformatics
  - Flood Risk Management





Global Change - Hydroinformatics - Planning



## Thank you for your attention





Global Change - Hydroinformatics - Planning