

# Theoretical overview : Disentangled uncertainty quantification for regression

Kevin PASINI et ALL

December 18, 2023

## Contents

<b>A Uncertainty in Regression task</b>	<b>1</b>
A.1 Detailed explanation of a disentangled uncertainty quantification formalism : . . . . .	3
<b>B Anomaly score :</b>	<b>5</b>
<b>C Additional mathematical justification :</b>	<b>6</b>
C.1 Notation . . . . .	6
C.2 Biases-variance trade-off and total uncertainty law: . . . . .	6
C.3 dE-Indicator development . . . . .	7

This document contains the theoretical sources related to the work on uncertainty quantification decomposed on ML/deep learning model for regression and anomaly detection.

## A Uncertainty in Regression task

Let us consider a modeling framework, in which random variables (denoted  $\varepsilon$ ) are linked to specific uncertainty sources. Here, time series forecasting (link to the type of data on which the work was undertaken) is treated as a regression problem based on time-dependent features. In this context, a model  $\hat{f}$  aims to predict the nominal behavior for variables of interest represented by univariate/multivariate time series  $Y = (y_1, \dots, y_t, \dots, y_T)$ . The forecast at time step  $t$  for the variable  $y_t$  will be based on a vector of observed variables  $\mathbf{x}_t$  composed of variables exogenous  $\mathbf{c}_t$  and lagged response  $\mathbf{Y}_t^{past} = (y_{t-lag}, \dots, y_{t-1})$  variables, as well as some latent variables  $\mathbf{h}_t$ :

$$y_t = \mathbf{f}(\mathbf{x}_t) + \varepsilon_t^u \quad \text{with; } \varepsilon_t^u \sim \mathcal{N}(0, \sigma_t^u(\mathbf{x}_t, \mathbf{h}_t)) \quad ; \quad \mathbf{x}_t = \{\mathbf{c}_t, \mathbf{Y}_t^{past}\}$$

with  $\mathbf{f}(\mathbf{x}_t)$  the average explainable signal, and  $\varepsilon_t^u$  a time-dependent Gaussian noise (local homogeneity assumption). The latter is associated with upstream irreducible variability, encompassing both intrinsic, measurement noises and pre-modeling noise arising from limits of the modeling scope (e.g. due to the influence of hidden variables  $\mathbf{h}_t$  that cannot be captured through lagged temporal variables  $\mathbf{Y}_t^{past}$ ).

The ML model  $\hat{f}_\theta$  aims to approximate the explainable part  $f$  of the target  $y$  using observed variables  $x$  from a training set  $D_\theta = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ , a subset of the dataset  $\mathcal{D}$ .  $\theta$  is the set of parameters obtained

using the training set  $D_\theta$ , over  $\Theta$  indicating the whole set of parameters linked to all subsets of the dataset  $\mathcal{D}$ . According to the bias-variance trade-off (Eq.1), we decompose all error sources between  $y_t$  and  $\hat{f}_\theta(\mathbf{x}_t)$ :

$$\begin{aligned} \mathbf{E}_\Theta [(y_t - \hat{f}_\theta(\mathbf{x}_t))^2] &= \mathbf{E}_\Theta [\hat{f}_\theta(\mathbf{x}_t) - f(\mathbf{x}_t)]^2 + \mathbf{E}_\Theta [(\mathbf{E}_\Theta [\hat{f}_\theta(\mathbf{x}_t)] - \hat{f}_\theta(\mathbf{x}_t))^2] + E_y [(y_t - f(\mathbf{x}_t))^2] \\ &= \underbrace{(f_\Theta^*(\mathbf{x}_t) - f(\mathbf{x}_t))^2}_{\text{Bias}} + \underbrace{\mathbf{E}_\Theta [(f_\Theta^*(\mathbf{x}_t) - \hat{f}_\theta(\mathbf{x}_t))^2]}_{\text{Variance}} + \underbrace{\sigma_t^u}_{\text{Intrinsic variability}}, \end{aligned} \quad (1)$$

with  $f_\Theta^* = \mathbf{E}_\Theta[\hat{f}_\theta(\mathbf{x}_t)]$ , the average function given the distribution  $\Theta$ . Among the three above-mentioned error sources, the *variance* can be explained by a noise  $\varepsilon_t^\theta$  which corresponds to the gap between the average function over  $\Theta$  and the ML model:  $\varepsilon_t^\theta = f_\Theta^*(\mathbf{x}_t) - \hat{f}_\theta(\mathbf{x}_t)$ . This epistemic noise is related to insufficient observations and could be reduced by gathering more data. The *bias* requires another random variable  $\varepsilon_t^\ominus$  linked to the gap between the average explainable signal and the average function over  $\Theta$ :  $\varepsilon_t^\ominus = f(\mathbf{x}_t) - f_\Theta^*(\mathbf{x}_t)$ . This noise, due to the modeling constraint over  $\Theta$ , is irreducible in the modeling scope. Finally, the *intrinsic variability* is related to the irreducible noise  $\varepsilon_t^u$  that appears upstream of the modeling scope. It quantifies a lower bound for the expected error in the test data with both infinite data and unconstrained modeling. To show the relation between the introduced random variables and the epistemic/aleatoric concepts, we inject them into the total uncertainty law [8] in Eq.11. After simplification allowed by strong independence and zeros-mean assumptions:

$$\begin{aligned} \text{With } y_t &= f_\theta(\mathbf{x}_t) + \varepsilon_t^\theta + \varepsilon_t^\ominus + \varepsilon_t^u \quad \text{and} \quad \sigma(y_t|x_t;\theta) = \sigma_\Theta [E_y(y_t|x_t;\theta)] + \mathbf{E}_\Theta [\sigma_y(y_t|x_t;\theta)] = \sigma_t^E + \sigma_t^A \\ \text{We obtain } \sigma_t^E &= \sigma_\Theta [\hat{f}_\theta(\mathbf{x}_t)] = \mathbf{E}_\Theta \left[ (\varepsilon_t^\theta)^2 \right] \quad , \quad \sigma_t^A = \sigma_y(\varepsilon_t^u) + \sigma_y(\varepsilon_t^\ominus) \end{aligned} \quad (2)$$

From these equations, we can see that the decomposition into epistemic and aleatoric components (denoted by  $E$  and  $A$  superscripts) requires the manipulation of the whole set of parameters  $\Theta$ . As expected, the epistemic part is essentially made up of the variance error caused by the sampling of the training set. However, the aleatoric part contains several quantities that are all irreducible in the modeling scope but may be associated with different sources: upstream modeling scope (intrinsic, measurement, and pre-modeling noise), and model constraints which also cause bias. When we move slightly outside the domain of validity of the assumptions (due to limited training data and approximate manipulation of  $\Theta$ ), the previous negligible terms can then induce blurs into the uncertainty decomposition.

**View of an unified dUQ framework:** The functional scheme of the proposed dUQ framework incorporating various UQ paradigms is shown in Fig. 2. It is based on a metamodel  $\mathcal{M}^\Theta$  that learns and manipulates diverse submodels ( $\hat{f}_\theta$ ) to combine their inferences. The learning phase aims to capture the explainable variability and estimate irreducible variability while exploring a diversity of submodel candidates  $\Theta$ . To ensure diversity and avoid submodel redundancy, a variability infusion mechanism (depending on the UQ paradigm) is needed during the learning phase. The estimated submodels produce, at the inference step, a local regressor  $\hat{f}_\theta(\mathbf{x}_t)$  and an estimation of aleatoric variability  $\hat{\sigma}^a(f_\theta(\mathbf{x}_t))$ . Furthermore, an epistemic variability  $\hat{\sigma}^e$  is produced by computing the variability of the submodel regression  $\hat{y}_t$  (using for example a Gaussian assumption). Finally, the metamodel provides a *risk-aware forecast* comprising three indicators:  $\bar{\mu}_t, \bar{\sigma}^a_t, \bar{\sigma}^e_t$  expressing forecast, aleatoric and epistemic indicators. As can be seen in Fig. 1, these indicators correspond to three independent axes on how the model perceives the data regarding forecast and sources of uncertainty. We can use them to design confidence intervals, error margins, or warnings highlighting a lack of model confidence.

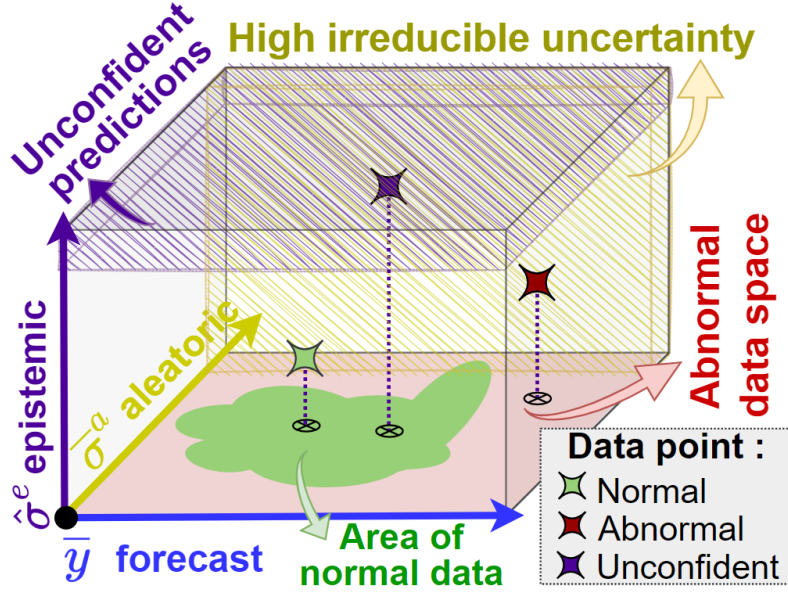


Figure 1: Theoretical *disentangled* Uncertainty Quantification (dUQ) indicators space

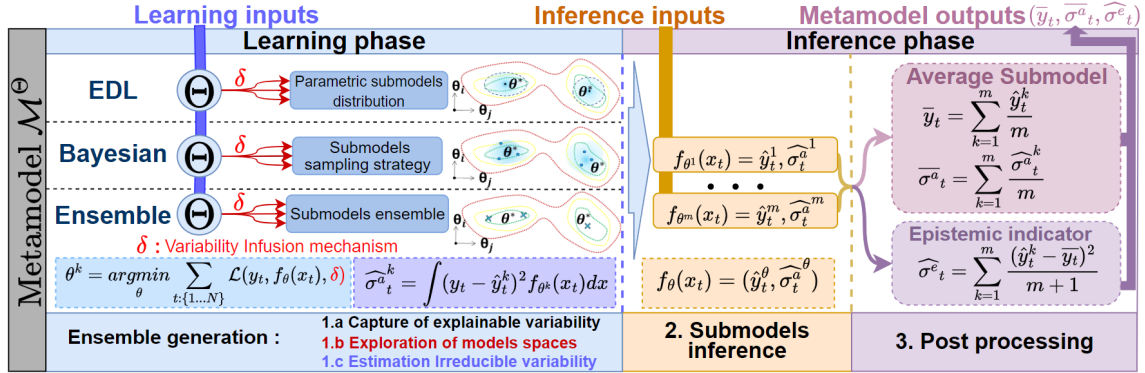


Figure 2: Illustration of a metamodel using Gaussian Aleatoric and Epistemic assumptions.

### A.1 Detailed explanation of a disentangled uncertainty quantification formalism :

**Submodels** are functions  $f = \{f^{\theta^1}, \dots, f^{\theta^m}\}$  parameterized by  $\theta$  that allow the *local* estimation of a decision function and its irreducible (Aleatoric) variability  $\varepsilon_t^A$  around the time stamp  $t$ . We have seen that this quantity explains the inherent noise in the data, which could not be explained during the training of the models but could be measured at each time stamp  $t$  for submodel  $k$ :

$$\varepsilon_t^A \sim \mathcal{N}(0, \sigma_t^a) : \sigma_t^a = \int (y_t - \mu_t)^2 f^{\theta^k}(\mathbf{x}_t) dy = \int (y_t - \hat{\mu}_t)^2 f^{\theta^k}(\mathbf{x}_t) dy = \widehat{\sigma}_t^a. \quad (3)$$

Such mechanisms take various forms in the literature (Bayesian model, Frequentist or Set-based approach) under different hypotheses (distribution assumption, moment estimation, quantiles, etc.). With a local Gaussian Aleatoric assumption justified by the application of central limit theorem on the observed variables, our submodels estimate the conditional probability distribution of the target variable from the observed variables.

$$\theta^k = \arg \min_{\theta} \sum_{t=1}^T \text{Loss}(y_t, \hat{\mu}_t^\theta, \widehat{\sigma}_t^{\alpha\theta}) \quad , \quad P_{\theta^k}(y_t|x_t) \sim \mathcal{N}(\hat{\mu}_t^k, \widehat{\sigma}_t^{\alpha k})$$

**The metamodel**  $\mathcal{M}^\theta$  is composed of a submodel ensemble (or submodel family). To specify a well-constructed metamodel, we have to introduce an abstract set  $\Omega$ : the family of reachable functions (that can be approximated by the type of selected model) and relevant submodels regarding a standard ML framework  $S$ . Then, the metamodel generates, manipulates, and enhances this ensemble of submodels through four following steps:

**1. The generation of the submodels ensemble** is realized by introducing variability during parameter optimization  $\theta$  through a variability infusion mechanism  $\delta$ , which takes different forms depending on the UQ-paradigm. The mechanisms of our four evaluated approaches can be interpreted as bagging/bootstrap (Ensemble), local minimum convergence (Deep Ensemble), binomial probabilistic weight (MCDP), or parameterization evidence-based (EDL). Mathematically, this is often achieved by disturbing the learning process (eq. 4) to generate submodel variants. To avoid having redundant submodels (lack of diversity) or coarse estimations (lack of accuracy), the metamodel has to be well-defined meaning:

$$\{\mathcal{M}^\theta = \mathcal{M}_{\{\theta^1, \dots, \theta^m\}} \text{ with } \theta^k \in \Omega \mid \theta^k = \arg \min_{\theta} \text{Loss}(Y, f^\theta(\mathbf{x}), \delta(Y, X, f^\theta))\} \quad \text{With} \quad (4)$$

- Considerable diversity of the submodels:  $\forall (k, k') \text{ with } k \neq k', \quad d(\theta^k, \theta^{k'}) > \epsilon$
- Acceptable accuracy and generalization capacity of submodels:  $\forall k \in [1, m], \quad f^{\theta^k} \in \Omega$

**2. Average model extrapolation** could be virtually computed through the average decision and the average aleatoric indicator of the submodels ensemble in a Gaussian aleatoric assumption.

$$\bar{\mu}_t = \frac{1}{m} \sum_{k=1}^m \hat{\mu}_t^k, \quad \bar{\sigma}_t^\alpha = \sqrt{\frac{1}{m} \sum_{k=1}^m (\widehat{\sigma}_t^{\alpha k})^2} \quad (5)$$

**3. Epistemic estimation**  $\varepsilon_t^E$  allows for local estimation of uncertainty through the metamodel submodels. Epistemic confidence can then be interpreted<sup>1</sup> as the likelihood of the *average submodel decision* ( $\bar{\mu}_t$ ) over the metamodel. Under the well-formed ensemble assumption for a regression task, local decision estimators  $f^{\theta^k}$  are centred with respect to the local mean  $\mu_t$  and have independent errors. Then, a Gaussian epistemic assumption<sup>2</sup> allows the use of the following unbiased empirical variance estimators.

$$\varepsilon_t^E \sim \mathcal{N}(0, \sigma_t^e): \quad \sigma_t^e \approx \widehat{\sigma}_t^e = \sqrt{\frac{1}{m+1} \sum_{k=1}^m (\hat{\mu}_t^k - \bar{\mu}_t)^2} \quad (6)$$

**4. Metamodel output** can express a decision and estimate its forecast uncertainty,  $\mathcal{M}_{\{\theta^k\}_{k=1}^m}(\mathbf{x}_t) \sim \mathcal{N}(\bar{\mu}_t, \widehat{\sigma}_t)$  corresponding to the combination of aleatoric and epistemic uncertainties  $\widehat{\sigma}_t^2 = \bar{\sigma}_t^{\alpha 2} + \widehat{\sigma}_t^{e 2}$ .

<sup>1</sup>More information in Section 3 of supplementary materials

<sup>2</sup>Thanks to unbiased (target-centred) models with independent errors

It should be noticed that there is a high correlation between the epistemic and aleatoric uncertainty indicators: the model error increases with high irreducible uncertainty. To better exploit local epistemic variability, we propose the dE-Indicator<sup>1</sup> corresponding to a Negative Ratio of Epistemic under Total Log-Likelihoods through Aleatoric and Epistemic Gaussian assumptions:  $I_t^e = -\ln(1 + \frac{\widehat{\sigma}_t^a}{\sigma_t^e})$

**Additional remarks** Two additional remarks should be made. First, the dUQ framework may be extended with other Aleatoric and Epistemic assumptions (distribution law, quantile, other kinds of moments). These extensions will lead to three main issues: submodel aleatoric estimation, submodel aggregation, and epistemic extraction from the ensemble. Secondly, this framework “contains” the standard *Machine Learning* (ML) approaches if they are formalized under the constant aleatoric assumption and total confidence in the unique submodel. It also “contains” other classical ML *Uncertainty Quantification* (UQ) focusing on Aleatoric or Epistemic if formalized under the assumption that neglecting the other part.

## B Anomaly score :

One standard way of characterizing abnormal behavior in a time series is to analyze the deviation, the difference between the observation and the behavior expected according to a model making a nominal prediction. It allows producing anomaly score based on prediction residuals. Then, it is possible to include consideration of an Uncertainty measure the score by creating a Z-score by normalizing the deviation based on the uncertainty to obtain a contextual anomaly score. In what follows, we propose a contextual anomaly score design including uncertainty, aleatorics and epistemics :

The design of the score anomaly ( $S_t$ ) is then based on three terms :

- The absolute deviation ( $r_t$ ) measures the normality of a sample.
- $\beta$ -level confidence interval threshold  $\widehat{\sigma}_t^\beta$  discriminate local statistical anomalies using local variability estimated by the model.
- $\alpha$ -level model risk penalization ( $pen_t^{E,\alpha}$ ) allows handling of predictions with low epistemic reliability.

$$S_t = \text{sign}(r_t) * \frac{|r_t| + pen_t^{E,\alpha}}{\widehat{\sigma}_t^\beta} : \begin{cases} r_t & = & (y_t - \widehat{y}_t) \\ \widehat{\sigma}_t^\beta & = & \beta * \sqrt{\widehat{\sigma}_t^A + \widehat{\sigma}_t^E} \\ pen_t^{E,\alpha} & = & \alpha * \widehat{\sigma}_t^E \end{cases} \quad (7)$$

$$S_t = \frac{(y_t - M(x_t)) \pm pen_t^{E,\alpha}}{\widehat{\sigma}_t^\beta} = \underbrace{\frac{e_t^a}{\widehat{\sigma}_t^\beta}}_{Ctx\ deviation} \pm \underbrace{\frac{\epsilon_t^A}{\widehat{\sigma}_t^\beta}}_{Data\ Noise} \pm \underbrace{\frac{B_t + \epsilon_t^E}{\widehat{\sigma}_t^\beta}}_{Model\ Noise} \pm \underbrace{\frac{pen_t^{E,\alpha}}{\widehat{\sigma}_t^\beta}}_{Model\ risk} \quad (8)$$

The following 3 assumptions allow us to make optimal use of the score:

- (i) Rare and significant anomaly :  $\forall t_i \text{ with } t_i \text{ is anom, } e_{t_i}^a \gg \widehat{\sigma}_{t_i}^\beta$
- (iI) good  $\sigma$ -estimator :  $\forall (t_i, t) \text{ with } t_i \text{ is anom, } \frac{e_{t_i}^a}{\widehat{\sigma}_{t_i}^\beta} \gg \frac{\epsilon_t^A + \epsilon_t^E}{\widehat{\sigma}_t^\beta}$  (9)
- (iii) good  $\mu$ -estimator :  $\widehat{\sigma}_t^\beta > B_t$  and  $B_t \Rightarrow 0$

Under the hypotheses described by formula 9, the model & data noises are negligible compared to the anomaly error in an abnormal context, while the noises are lower than the detection threshold in the nominal situation due to normalization of the variance. Therefore, the anomaly score  $S_t$  discriminates ( $|S_t| > 1$ ) statistical anomalies from a contextual confidence interval defined for a confidence threshold  $\beta$  with particular consideration for the risk of prediction unreliability according to a specified critical value  $\alpha$  (**misprediction risk aversion**).

## C Additional mathematical justification :

### C.1 Notation

NOTATIONS	
Basic term	
$t$	index associated with a time step of a time series
$k$	index associated to a member (submodel) of a metamodel
$N$	number of samples (time series)
$m$	number of submodels in a metamodel
$\hat{a}$	Estimation of the quantity $a$
$\bar{a}$	Average of the quantity $a$ (often used on a set of estimations)
ML Modeling variable	
$y_t$	the value of the series $y$ (ground truth) at time step $t$
$f$	function expressing the explainable part of the series $y$
$c_t, h_t$	exogenous variables (observed, hidden) impacting the time series $y_t$
$Y_t^{past}$	lagged endogenous variables from series $y$ used for forecasting of future values
$\mathbf{x}_t$	the input data at time stamps $t$ composed in our framework by $(c_t, Y_t^{past})$
$D, D_\theta$	dataset composed of pairs $(X_1, y_1), \dots, (X_N, y_N)$ , and subset of dataset used to train a model
$\Theta$	Ensemble of parameters set
$\theta, \theta^i, \theta^*$	parameters set (resp one, the $i^{th}$ , the "optimal") from $\Theta$
$\hat{f}_\theta(\mathbf{x}_t) = \hat{y}_t$	a ML model (function approximation) parametrized by $\theta$ with $x$ as inputs and providing $\hat{y}$ as output
UQ modeling term	
$\delta$	abstract diversity infusion mechanism allowing the exploration of $\Theta$ space (ex bootstrap).
$\mathcal{M}, \mathcal{M}^c, \mathcal{M}^d$	metamodel, control metamodel, degraded metamodel impacted by a variability injection
$\varepsilon^u, \varepsilon^\theta, \varepsilon^\beta$	Noise (uncertainty) link respectively to upstream sources, modeling constraint, and epistemic issues
$\sigma_t^A, \sigma_t^E$	Variance (Aleatoric, Epistemic), at step $t$
model and meta-model output	
$\hat{f}_\theta(\mathbf{x}_t) = (\hat{y}_t^{\theta}, \hat{\sigma}_t^{A,\theta})$	Outputs (Forecast and aleatoric estimation) corresponding to submodel $\theta$ at time step $t$
$M(\mathbf{x}_t) = (\bar{y}, \bar{\sigma}_t^A, \bar{\sigma}_t^E)$	Outputs (Forecast, aleatoric, epistemic estimation) corresponding to metamodel $M$ at time step $t$
Experimental modeling term	
$I_t^e$	disentangled Epistemic indicator (dE-Indicator)
$X_n, X_a$	nominal query without changes, altered queries that are degraded by variability injection

### C.2 Biases-variance trade-off and total uncertainty law:

Taking up the formalization introduced in the paper, we introduced three aleatoric variable that can each be associated with a specific source of errors in the Biases-variance trade off.

- $\varepsilon_t^u = y_t - f(\mathbf{x}_t)$  the noise upstream the modeling scope associated with intrinsic irreducible variability errors, that corresponds to the difference between the series  $Y$ , and its explainable part  $f$  according to the modeling scope  $S$ .

- $\varepsilon_t^\Theta = f(\mathbf{x}_t) - f_\Theta^*(\mathbf{x}_t)$  the modeling-constraint noise, related to bias errors due to the exploration space  $\Theta$ , that correspond to the difference between the explainable part  $f$ – of data and the optimal reachable function  $f_\Theta^*$  in  $\Theta$ .
- $\varepsilon_t^\theta = f_\Theta^*(\mathbf{x}_t) - \widehat{f}_\theta(\mathbf{x}_t)$  the epistemic noise, related to variance errors, that correspond to the difference between the optimal reachable function  $f_\Theta^*$  in  $\Theta$ , and the obtained model  $\widehat{f}_\theta$  which is biased by limited dataset  $D_\theta$ .

This modeling frame provides us with an artificial way to decompose the series  $y$  according to a model, and its different sources of errors :

$$\text{With: } y_t = f(\mathbf{x}_t) + \varepsilon_t^u = f_\Theta^*(\mathbf{x}_t) + \varepsilon_t^\Theta + \varepsilon_t^u = \widehat{f}_\theta(\mathbf{x}_t) + \varepsilon_t^\theta + \varepsilon_t^\Theta + \varepsilon_t^u \quad (10)$$

This modeling frame provides us with an artificial way to decompose the series  $y$  according to a model, and its different sources of errors. It can then be rejected in the total law of uncertainty decomposition to analyze how this term impacts an aleatoric versus epistemic decomposition :

$$\begin{aligned} \text{With } \sigma(y_t|x_t;\theta) &= \sigma_\Theta [E_y(y_t|x_t;\theta)] + E_\Theta [\sigma_y(y_t|x_t;\theta)] = \sigma_t^E + \sigma_t^A \\ \sigma_t^E &\stackrel{\text{if ind.}}{=} \underbrace{\sigma_\Theta [E_y[\widehat{f}_\Theta^{\theta k}(\mathbf{x}_t)]]}_{E_\Theta [(f_\Theta^*(\mathbf{x}_t) - \widehat{f}_\theta(\mathbf{x}_t))^2]} + \underbrace{\sigma_\Theta [E_y[\varepsilon_t^\Theta]]}_{0^*} + \underbrace{\sigma_\Theta [E_y[\varepsilon_t^u]]}_{0^*} + \underbrace{\sigma_\Theta [E_y[\varepsilon_t^\theta]]}_{0^*} \\ \sigma_t^A &\stackrel{\text{if ind.}}{=} \underbrace{E_\Theta [\sigma_y(f_\Theta^{\theta k}(\mathbf{x}_t))]}_0 + \underbrace{\sigma_y(\varepsilon_t^\Theta)}_{\text{modeling}} + \underbrace{\sigma_y(\varepsilon_t^u)}_{\text{upstream}} + \underbrace{E_\Theta [E_y(\varepsilon_t^\theta - E_y[\varepsilon_t^\theta])^2]}_{0^*} \end{aligned} \quad (11)$$

Two types of assumption allow making simplification. We can consider that the random variables are independent (ind.), obtaining independent terms, then consider that the different random variables are centered (for different reasons), which makes the terms negligible.

### C.3 dE-Indicator development

By manipulating dUQ framework under assumptions about aleatoric and epistemic uncertainty, we can design predictive indicators from metamodel likelihood using the assumption A1: the submodels ensemble of metamodel is well-formed, which means rich in variability  $\forall (k, k')$  with  $k \neq k' \quad d(\theta^k, \theta^{k'}) > c$  and composed of submodels with good predictions and with few overfitting  $\forall k, f_{\theta^k} \in \Theta$ . Starting from the normal log-likelihood using total variance  $\sigma^{tot}$  :

$$\begin{aligned} \ln(L(\bar{y}_t, \widehat{\sigma}^{tot}_t; \Theta, y_t)) &= \ln(P(y_t; \bar{y}_t, \widehat{\sigma}^{tot}_t | \Theta)) \\ &\approx \ln(P(y_t; \bar{y}_t, \widehat{\sigma}^{tot}_t | \{\theta^k\}_{k \in [1, m]})) \\ &= cst - \ln(\widehat{\sigma}^{tot}_t) - \frac{(y_t - \bar{y}_t)^2}{2 * \widehat{\sigma}^{tot}_t} \end{aligned} \quad \left. \begin{array}{l} \text{A1 submodel-ensemble approx} \\ \text{NLL developpement} \end{array} \right\} \quad (12)$$

A local likelihood of the metamodel on the training set can be approximated under:

- the normal-aleatoric assumption  $y_t \sim \mathcal{N}(\bar{y}_t, \sigma^a)$

- the normal-epistemic assumptions  $\mu_t^* \sim \mathcal{N}(\bar{y}_t, \bar{\sigma}^e)$
- the A2 assumption of a local unbiased capture based on the ability of the submodels to perform local estimation in a standard ML-framework.

We can derive the A2 assumption of a "local unbiased capture performed by submodels" by combining a standard ML-framework and normal-aleatoric assumption.

$$\text{A2: Local neighbor observation approx: } y_t \approx \int \sum_k f_{\theta^k}(\mathbf{x}_t) dx \approx \int y_t dy$$

$$\begin{aligned} \ln(L(\bar{y}_t, \widehat{\sigma}^{tot}_t; \Theta, y_t)) &\approx cst - \ln(\widehat{\sigma}^{tot}_t) - \frac{(y_t - \bar{y}_t)^2}{2 * \widehat{\sigma}^{tot}_t} \\ &\approx cst - \ln(\widehat{\sigma}^{tot}_t) - \frac{(\int y_t dy - \bar{y}_t)^2}{2 * \widehat{\sigma}^{tot}_t} \\ &\approx cst - \ln(\widehat{\sigma}^{tot}_t) + 0 \\ &\approx cst - \ln(\bar{\sigma}^a_t + \widehat{\sigma}^e_t) + 0 \end{aligned} \quad \left. \begin{array}{l} \text{Formula 5 + } y_t \sim A2 \\ \bar{y}_t: \text{Local MLE by NLL loss minimization} \end{array} \right\} \quad (13)$$

$$L(\bar{y}_t, \widehat{\sigma}^{tot}_t | \Theta) \propto \frac{1}{\bar{\sigma}^a_t + \widehat{\sigma}^e_t}$$

A local epistemic likelihood of the "average-submodel" mean among the reachable and relevant models  $\Theta$  can also be obtained under the A3 assumption of "local unbiased mean modelling approximation" stemming from the combination of the A1 assumption and normal-epistemic assumptions.

$$\text{A3 Local mean modelling approx: } \mu_t^* \approx \int_{\Theta} \frac{\hat{y}_t^k}{m} \approx \sum_{k=1}^m \frac{\hat{y}_t^k}{m}$$

$$\begin{aligned} \ln(L(\bar{y}_t, \widehat{\sigma}^e_t | \Theta, \mu_t^*)) &= cst - \ln(\widehat{\sigma}^e_t) - \frac{(\mu_t^* - \bar{y}_t)^2}{2 * \widehat{\sigma}^e_t} \\ &\approx cst - \ln(\widehat{\sigma}^e_t) - \frac{\left(\sum_{k=1}^m \frac{\hat{y}_t^k}{m} - \bar{y}_t\right)^2}{2 * \widehat{\sigma}^e_t} \\ &= cst - \ln(\widehat{\sigma}^e_t) - 0 \end{aligned} \quad \left. \begin{array}{l} \text{Formula 5 + } \mu_t^* \sim A3 \\ \bar{y}_t: \text{local MLE} \end{array} \right\} \quad (14)$$

$$L(\bar{y}_t, \widehat{\sigma}^e_t | \Theta) \propto \frac{1}{\widehat{\sigma}^e_t}$$

Finally, we are more specifically interested in the fluctuation of the local epistemic likelihood of the decision regarding local forecasting uncertainty to be more "invariant" to the magnitude of local forecasting uncertainty. Therefore, we propose to analyze the local epistemic likelihood indicator (Formula 14) normalized by the local total likelihood indicator (Formula 13).

$$\begin{aligned} \frac{L^{epi}(\bar{y}_t, \widehat{\sigma}^e_t | \Theta)}{L^{tot}(\bar{y}_t, \widehat{\sigma}^{tot}_t | \Theta)} &\propto \frac{\bar{\sigma}^a_t + \widehat{\sigma}^e_t}{\widehat{\sigma}^e_t} \\ \ln\left(\frac{L^{epi}}{L^{tot}}\right) &\propto \ln(\bar{\sigma}^a_t + \widehat{\sigma}^e_t) - \ln(\widehat{\sigma}^e_t) \\ &\propto \ln\left(1 + \frac{\bar{\sigma}^a_t}{\widehat{\sigma}^e_t}\right) \end{aligned} \quad (15)$$



## References

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: a system for large-scale machine learning. In *OsdI*, volume 16, pp. 265–283. Savannah, GA, USA, 2016.
- [2] Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad Ghavamzadeh, Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U Rajendra Acharya, et al. A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information Fusion*, 2021.
- [3] Alexander Amini, Wilko Schwarting, Ava Soleimany, and Daniela Rus. Deep evidential regression. *Advances in Neural Information Processing Systems*, 33:14927–14937, 2020.
- [4] Arsenii Ashukha, Alexander Lyzhov, Dmitry Molchanov, and Dmitry Vetrov. Pitfalls of in-domain uncertainty estimation and ensembling in deep learning. *arXiv preprint arXiv:2002.06470*, 2020.
- [5] Sheraz Aslam, Herodotos Herodotou, Syed Muhammad Mohsin, Nadeem Javaid, Nouman Ashraf, and Shahzad Aslam. A survey on deep learning methods for power load and renewable energy forecasting in smart microgrids. *Renewable and Sustainable Energy Reviews*, 144:110992, 2021.
- [6] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In *International conference on machine learning*, pp. 1613–1622. PMLR, 2015.
- [7] Ian C. Covert, Scott Lundberg, and Su-In Lee. Explaining by removing: a unified framework for model explanation. *The Journal of Machine Learning Research*, 22(1):209:9477–209:9566, January 2021. ISSN 1532-4435.
- [8] Stefan Depeweg, Jose-Miguel Hernandez-Lobato, Finale Doshi-Velez, and Steffen Udluft. Decomposition of uncertainty in bayesian deep learning for efficient and risk-sensitive learning. In *International Conference on Machine Learning*, pp. 1184–1193. PMLR, 2018.
- [9] Armen Der Kiureghian and Ove Ditlevsen. Aleatory or epistemic? does it matter? *Structural safety*, 31(2):105–112, 2009.
- [10] Timothy Dozat. Incorporating nesterov momentum into adam. 2016.
- [11] Romain Egele, Romit Maulik, Krishnan Raghavan, Bethany Lusch, Isabelle Guyon, and Prasanna Balaprakash. Autodeuq: Automated deep ensemble with uncertainty quantification. In *2022 26th International Conference on Pattern Recognition (ICPR)*, pp. 1908–1914. IEEE, 2022.
- [12] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pp. 1050–1059. PMLR, 2016.
- [13] Jakob Gawlikowski, Cedrique Rovile Njieutcheu Tassi, Mohsin Ali, Jongseok Lee, Matthias Humt, Jianxiang Feng, Anna Kruspe, Rudolph Triebel, Peter Jung, Ribana Roscher, et al. A survey of uncertainty in deep neural networks. *arXiv preprint arXiv:2107.03342*, 2021.
- [14] Alex Graves. Practical variational inference for neural networks. *Advances in neural information processing systems*, 24, 2011.

- [15] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pp. 1321–1330. PMLR, 2017.
- [16] Giles Hooker, Lucas Mentch, and Siyu Zhou. Unrestricted permutation forces extrapolation: variable importance requires at least one more model, or there is no free variable importance. *Statistics and Computing*, 31(6):82, October 2021. ISSN 1573-1375. doi: 10.1007/s11222-021-10057-z. URL <https://doi.org/10.1007/s11222-021-10057-z>.
- [17] Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods. *Machine Learning*, 110(3):457–506, 2021.
- [18] Thanasis Kotsiopoulos, Panagiotis Sarigiannidis, Dimosthenis Ioannidis, and Dimitrios Tzovaras. Machine learning and deep learning in smart manufacturing: The smart grid paradigm. *Computer Science Review*, 40:100341, 2021.
- [19] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.
- [20] Chan-Chiao Lin, Huei Peng, Soonil Jeon, and Jang Moo Lee. Control of a hybrid electric truck based on driving pattern recognition. In *Proceedings of the 2002 advanced vehicle control conference, Hiroshima, Japan*, pp. 9–13, 2002.
- [21] Jeremiah Liu, John Paisley, Marianthi-Anna Kioumourtzoglou, and Brent Coull. Accurate uncertainty estimation and decomposition in ensemble learning. *Advances in neural information processing systems*, 32, 2019.
- [22] David John Cameron Mackay. *Bayesian methods for adaptive models*. California Institute of Technology, 1992.
- [23] Nis Meinert and Alexander Lavin. Multivariate deep evidential regression. *arXiv preprint arXiv:2104.06135*, 2021.
- [24] Nicolai Meinshausen and Greg Ridgeway. Quantile regression forests. *Journal of machine learning research*, 7(6), 2006.
- [25] Mouhcine Mendil, Marc Mossina, Lucaand Nabhan, and Kevin Pasini. Robust gas demand forecasting with conformal prediction. In *Conformal and Probabilistic Prediction with Applications*, pp. 169–187. PMLR, 2022.
- [26] Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.
- [27] Giovanna Nicora, Miguel Rios, Ameen Abu-Hanna, and Riccardo Bellazzi. Evaluating pointwise reliability of machine learning prediction. *Journal of Biomedical Informatics*, 127:103996, 2022.
- [28] Dongpin Oh and Bonggun Shin. Improving evidential deep learning via multi-task learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 7895–7903, 2022.
- [29] Kevin PASINI. *Forecast and anomaly detection on time series with dynamic context : Application to the mining of transit ridership data*. PhD thesis, Université Gustave Eiffel, 2021.

- [30] Kevin Pasini, Mostepha Khouadjia, Allou Samé, Martin Trépanier, and Latifa Oukhellou. Contextual anomaly detection on time series: a case study of metro ridership analysis. *Neural Computing and Applications*, 34(2):1483–1507, 2022.
- [31] Tim Pearce, Mohamed Zaki, Alexandra Brintrup, N Anastassacos, and A Neely. Uncertainty in neural networks: Bayesian ensembling. *stat*, 1050:12, 2018.
- [32] Murat Sensoy, Lance Kaplan, and Melih Kandemir. Evidential deep learning to quantify classification uncertainty. *Advances in neural information processing systems*, 31, 2018.
- [33] Yuan Wang, Dongxiang Zhang, Ying Liu, Bo Dai, and Loo Hay Lee. Enhancing transportation systems via deep learning: A survey. *Transportation research part C: emerging technologies*, 99:144–163, 2019.
- [34] Florian Wenzel, Jasper Snoek, Dustin Tran, and Rodolphe Jenatton. Hyperparameter ensembles for robustness and uncertainty quantification. *Advances in Neural Information Processing Systems*, 33: 6514–6527, 2020.
- [35] Lingxue Zhu and Nikolay Laptev. Deep and confident prediction for time series at uber. In *2017 IEEE International Conference on Data Mining Workshops (ICDMW)*, pp. 103–110. IEEE, 2017.