

Evidential Machine Learning

Thierry Denœux

Université de technologie de Compiègne, Compiègne, France
Institut Universitaire de France, Paris, France

<https://www.hds.utc.fr/~tdenoeux>

7th School on Belief Functions and their Applications
Granada, Spain, October 23, 2025

Uncertainty quantification in machine learning

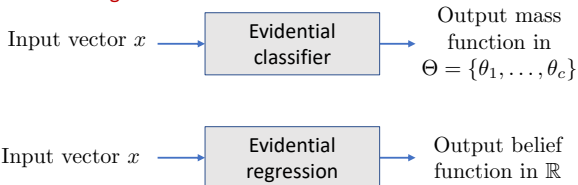
- Uncertainty quantification (UQ): a topical problem in machine learning
 - ▶ Prediction uncertainty (supervised learning)
 - ▶ Cluster-membership uncertainty (clustering)
- Main sources:
 - ▶ Randomness (aleatory uncertainty)
 - ▶ Lack of knowledge (epistemic uncertainty)
- Theoretical frameworks:
 - ▶ Frequentist statistical framework (prediction regions)
 - ▶ Bayesian (additive probabilities)
 - ▶ Imprecise probabilities (credal sets, lower previsions)
 - ▶ Dempster-Shafer theory (random sets, belief functions)
 - ▶ Generalised evidence theory (random fuzzy sets, belief functions)

Evidential Machine Learning

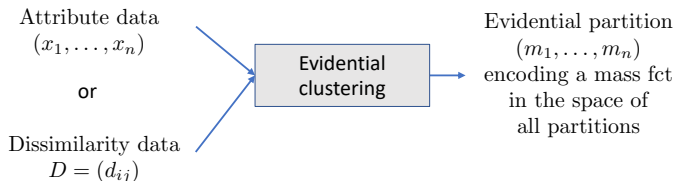
- Arguments for (generalised) evidence theory in ML:
 - ▶ Allows for the representation of aleatory and epistemic uncertainties
 - ▶ Does not require a precise prior as Bayesian inference (but yields the same results as Bayesian inference if a prior is provided)
 - ▶ Well-suited for information fusion (multimodal/multisensor classification, classifier/clusterer combination, etc.)
- Evidential Machine Learning is a branch of ML based on (standard or generalised) evidential reasoning. It quantifies uncertainty using belief functions.

Overview of Evidential ML

Supervised learning



Unsupervised learning



Distance-based vs. likelihood-based approaches

- In this talk, we focus on **supervised learning**.
- Two main approaches:
 - 1 The **distance-based** approach computes distances between input vectors and learning vectors or prototypes, seen as pieces of evidence.
 - Classification: evidential k-NN rule, evidential neural network
 - Regression: ENNreg model
 - 2 The **likelihood-based** approach assumes a parametric statistical model. The relative likelihood function is considered as a possibility distribution, which allows one to construct a **predictive random fuzzy set**.

Outline

- 1 Distance-based approach
 - Neural network model for classification
 - Neural network model for regression
- 2 Likelihood-based approach
 - Evidential Likelihood-based inference
 - Application to machine learning

Outline

- 1 Distance-based approach
 - Neural network model for classification
 - Neural network model for regression
- 2 Likelihood-based approach
 - Evidential Likelihood-based inference
 - Application to machine learning

Distance-based approach for classification

- Two of the first papers applying DS theory to classification:



T. Denœux

A k-nearest neighbor classification rule based on Dempster-Shafer theory
IEEE Transactions on SMC, 25(05):804–813, 1995

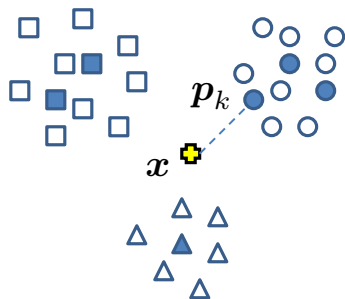


T. Denœux

A neural network classifier based on Dempster-Shafer theory
IEEE transactions on SMC A, 30(2):131–150, 2000

- I will briefly recall the evidential neural network model and describe some recent developments.

Evidential neural network classifier



- The learning set is summarized by K **prototypes**
- Each prototype p_k has **membership degree** $u_{kl} \geq 0$ to each class ω_l , with $\sum_{l=1}^c u_{kl} = 1$
- Each prototype p_k is a **piece of evidence** about the class of x ; its **reliability decreases with its distance to x**

Propagation equations

- Frame: $\Theta = \{\theta_1, \dots, \theta_c\}$ for c -class classification.
- Similarity between input \mathbf{x} and prototype \mathbf{p}_k :

$$s_k(\mathbf{x}) = \exp(-\gamma_k \|\mathbf{x} - \mathbf{p}_k\|^2), \quad \gamma_k \geq 0$$

- Mass function induced by prototype \mathbf{p}_k :

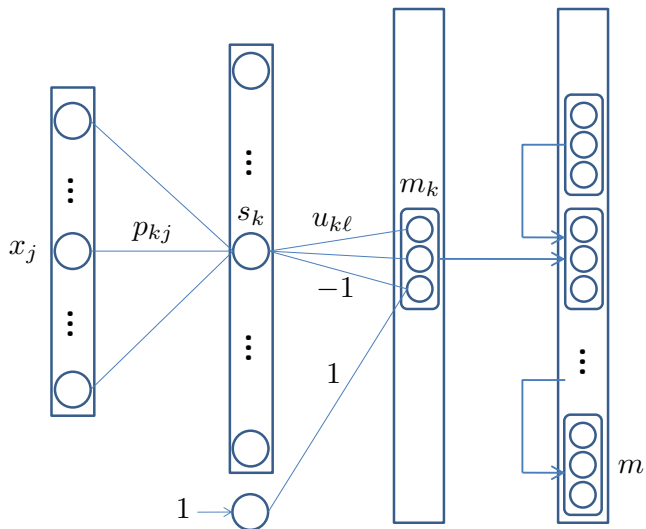
$$\begin{aligned} m_k(\{\theta_\ell\}) &= \alpha_k u_{k\ell} s_k(\mathbf{x}), \quad \ell = 1, \dots, c \\ m_k(\Theta) &= 1 - \alpha_k s_k(\mathbf{x}) \end{aligned}$$

with $\alpha_k \in [0, 1]$.

- Property: $\lim_{s_k \rightarrow 0} m_k(\Theta) = 1$.
- Combination of the evidence from the K prototypes:

$$m = \bigoplus_{k=1}^K m_k$$

Neural network implementation



Learning

- The parameters are the
 - ▶ The prototypes \mathbf{p}_k , $i = k, \dots, K$ (Kp parameters)
 - ▶ The membership degrees $u_{k\ell}$, $k = 1, \dots, K$, $\ell = 1, \dots, c$ (Kc parameters)
 - ▶ The α_k and γ_k , $k = 1, \dots, K$ ($2K$ parameters).
- The parameter vector Ψ can be estimated by minimising a **loss function** such as

$$C_\lambda(\Psi) := \underbrace{\sum_{i=1}^n \sum_{\ell=1}^c (pl_{i\ell} - y_{i\ell})^2}_{\text{error}} + \lambda \underbrace{\sum_{k=1}^K \alpha_k}_{\text{regularization}}$$

where $pl_{i\ell}$ is the output plausibility of class θ_ℓ for instance i , $y_{i\ell} = I(y_i = \theta_\ell)$, and λ is a regularisation coefficient (hyperparameter).

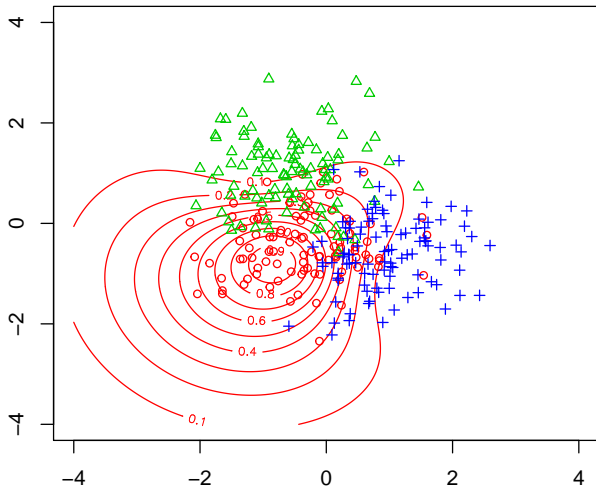
- λ can be optimised by cross-validation.

Implementations

- Matlab: `http://www.hds.utc.fr/~tdenoeux/software/belief_NN/belief_NN.zip`
- R package `evclass`, available at `https://cran.r-project.org/web/packages/evclass`

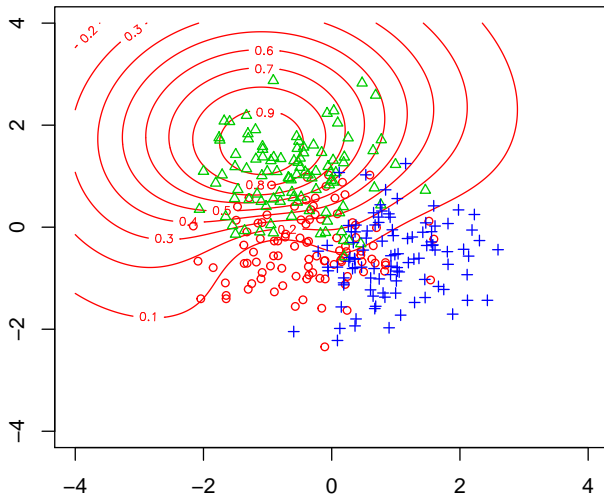
Example

Mass on $\{\theta_1\}$



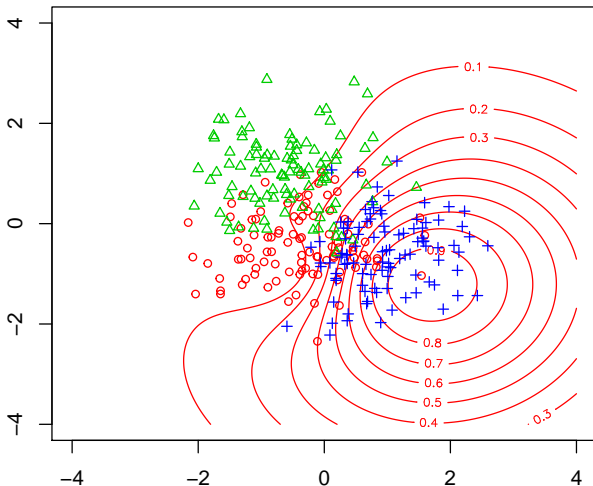
Example

Mass on $\{\theta_2\}$



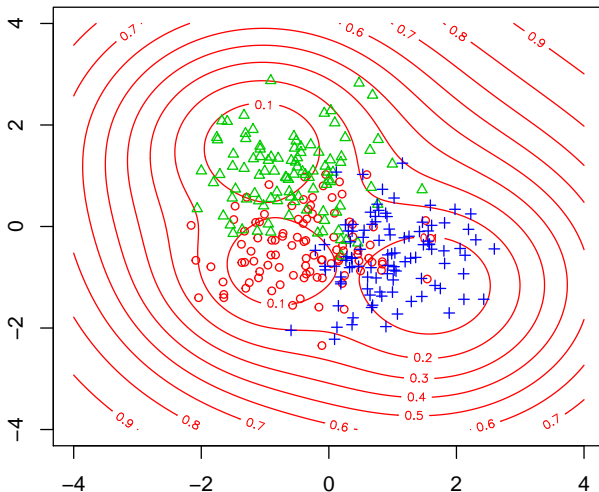
Example

Mass on $\{\theta_3\}$



Example

Mass on Θ



Decision

- The focal sets of the combined mass function m are the singletons $\{\theta_k\}$, $k = 1, \dots, c$ and Θ .
- For decision analysis, we need to define:
 - ▶ A set \mathcal{O} of outcomes (correct decision, error, no classification, etc)
 - ▶ A set of $\mathcal{G} = \{g_1, \dots, g_q\}$ of acts, where $g_k : \Theta \rightarrow \mathcal{O}$ (classification in class θ_k , partial classification in set $A \subset \Theta$, rejection, etc.)
 - ▶ A utility function $u : \mathcal{O} \rightarrow \mathbb{R}$
- Different decision rules for precise or partial classification can be defined, see



T. Denœux

Decision-Making with Belief Functions: a Review

International Journal of Approximate Reasoning, 109:87–110, 2019



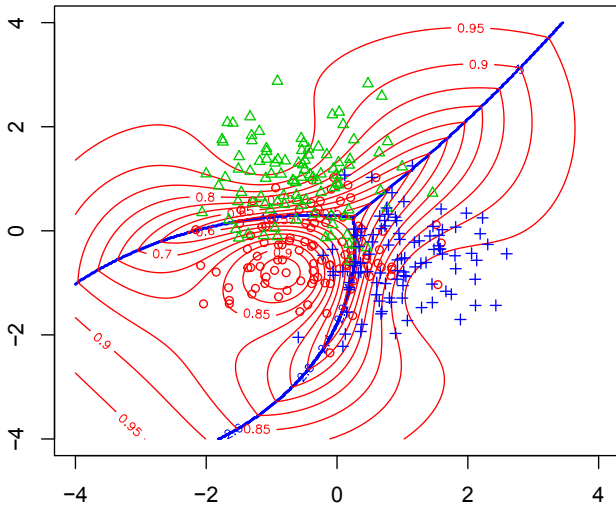
L. Ma and T. Denœux

Partial Classification in the Belief Function Framework

Knowledge-Based Systems, 214:106742, 2021

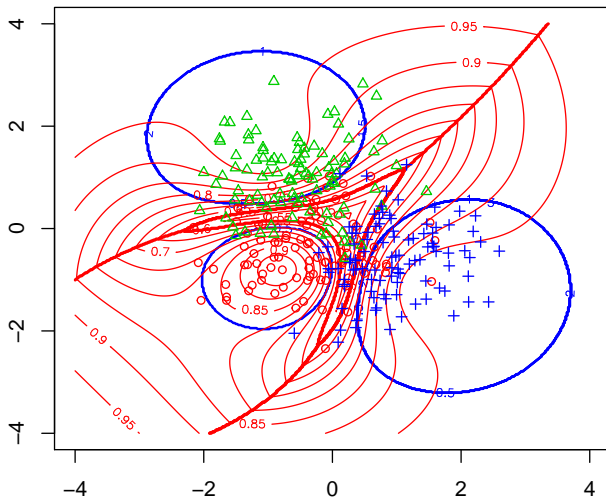
Example

Decision boundaries, optimistic/pessimistic decision rules (no rejection)



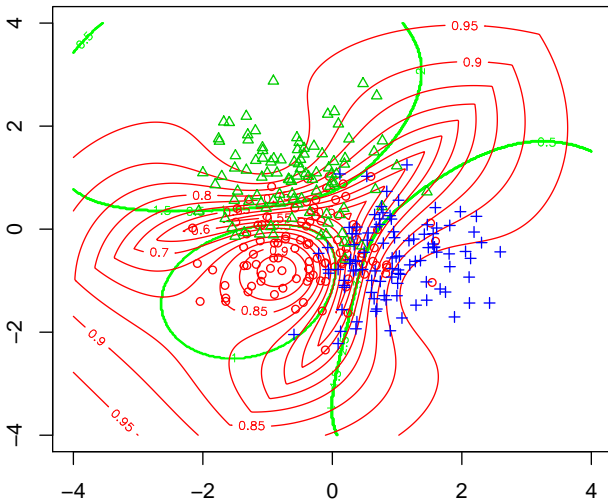
Example

Decision boundaries, optimistic/pessimistic decision rules (with rejection)

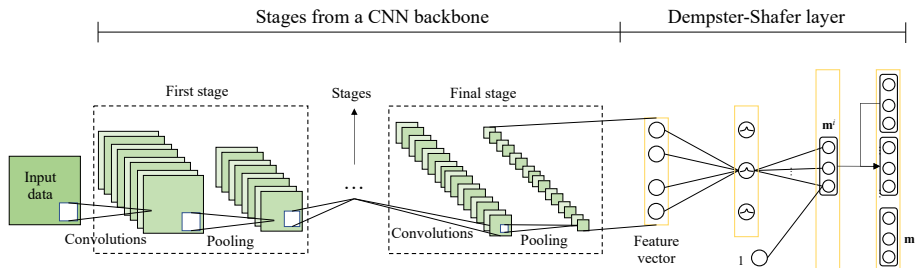


Example

Decision boundaries, Hurwicz criterion, $\rho = 0.5$ (with rejection)



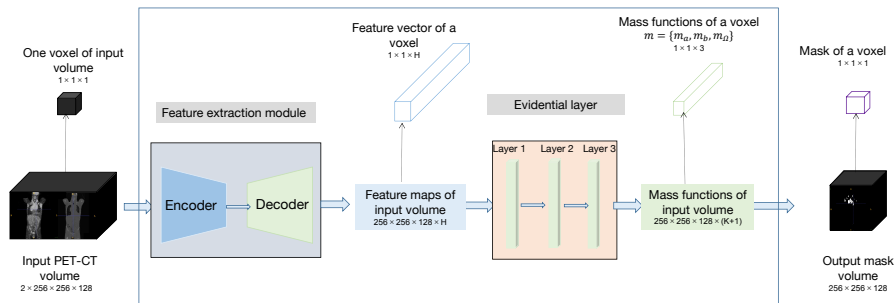
Deep evidential classifier



Z. Tong, Ph. Xu and T. Denœux

An evidential classifier based on Dempster-Shafer theory and deep learning
Neurocomputing 450:275–293, 2021

Application to medical image segmentation

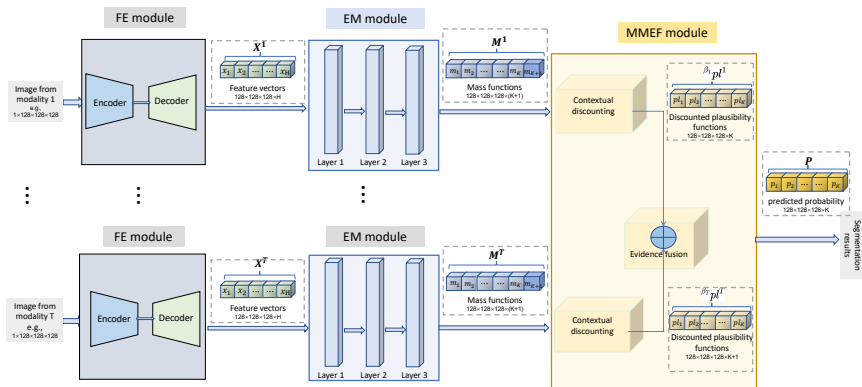


 L. Huang, S. Ruan, P. Decazes and T. Denœux

Lymphoma segmentation from 3D PET-CT images using a deep evidential network

International Journal of Approximate Reasoning 149:39-60, 2022

Application to multimodal medical image segmentation



L. Huang, S. Ruan, P. Decazes and T. Denœux

Deep evidential fusion with uncertainty quantification and reliability learning for multimodal medical image segmentation

Information Fusion, 113:102648, 2025

Outline

1 Distance-based approach


- Neural network model for classification
- Neural network model for regression

2 Likelihood-based approach

- Evidential Likelihood-based inference
- Application to machine learning

The ENNreg model

- We now consider a **regression problem**: the task is to predict a continuous random response Y from p input variables $\mathbf{X} = (X_1, \dots, X_p)$, based on a learning set $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$.
- ENNreg is a **neural network model**¹ that quantifies uncertainty about the response Y given input vector $\mathbf{X} = \mathbf{x}$ by a **GRFN** $\tilde{Y}(\mathbf{x})$ with associated belief function $Bel_{\tilde{Y}(\mathbf{x})}$.
- As the ENN model for classification, ENNreg is based on **prototypes**. The prototypes are treated as **independent pieces of evidence** about the response, combined by the product-intersection rule.

¹T. Denœux. Quantifying Prediction Uncertainty in Regression using Random Fuzzy Sets: the ENNreg model. *IEEE Transactions on Fuzzy Systems*, 31(10):3690–3699, 2023. 

Propagation equations (1/2)

- Let $\mathbf{p}_1, \dots, \mathbf{p}_K$ denote K prototypes in the p -dimensional input space.
- **Similarity** between input vector \mathbf{x} and prototype \mathbf{p}_k :

$$s_k(\mathbf{x}) = \exp(-\gamma_k^2 \|\mathbf{x} - \mathbf{p}_k\|^2)$$

where $\gamma_k > 0$ is a scale parameter.

- The **evidence from prototype \mathbf{p}_k** is represented by a GRFN

$$\tilde{Y}_k(\mathbf{x}) \sim \tilde{N}(\mu_k(\mathbf{x}), \sigma_k^2, s_k(\mathbf{x})h_k)$$

where σ_k^2 and h_k are variance and precision parameters, and

$$\mu_k(\mathbf{x}) = \boldsymbol{\beta}_k^T \mathbf{x} + \beta_{k0} \quad \text{with} \quad \boldsymbol{\beta}_k \in \mathbb{R}^p, \beta_{k0} \in \mathbb{R}$$

Propagation equations (2/2)

- The output $\tilde{Y}(\mathbf{x})$ for input \mathbf{x} is computed as

$$\tilde{Y}(\mathbf{x}) = \tilde{Y}_1(\mathbf{x}) \tilde{\oplus} \dots \tilde{\oplus} \tilde{Y}_K(\mathbf{x})$$

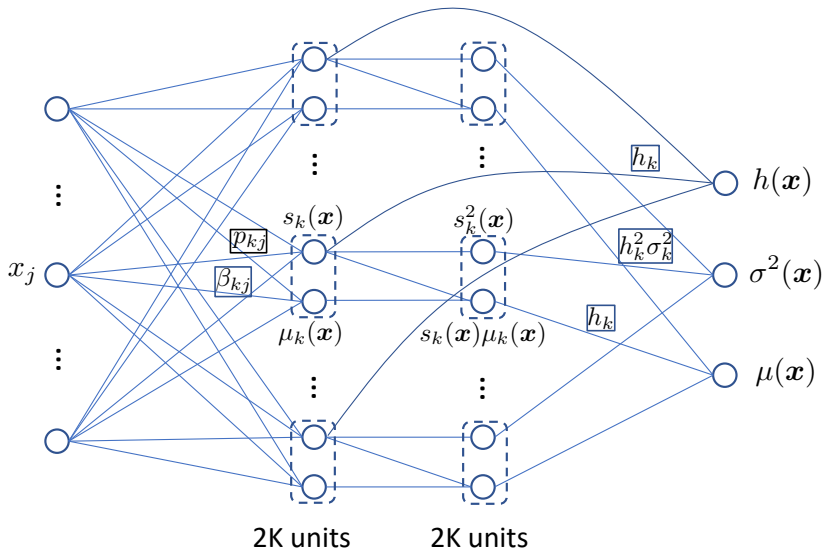
where $\tilde{\oplus}$ is an approximation of the product-intersection rule (assuming $\tilde{\mu} \approx \mu$ and $\tilde{\Sigma} \approx \Sigma$).

- We have $\tilde{Y}(\mathbf{x}) \sim \tilde{N}(\mu(\mathbf{x}), \sigma^2(\mathbf{x}), h(\mathbf{x}))$, with

$$\mu(\mathbf{x}) = \frac{\sum_{k=1}^K s_k(\mathbf{x}) h_k \mu_k(\mathbf{x})}{\sum_{k=1}^K s_k(\mathbf{x}) h_k}$$

$$\sigma^2(\mathbf{x}) = \frac{\sum_{k=1}^K s_k^2(\mathbf{x}) h_k^2 \sigma_k^2}{\left(\sum_{k=1}^K s_k(\mathbf{x}) h_k\right)^2} \quad \text{and} \quad h(\mathbf{x}) = \sum_{k=1}^K s_k(\mathbf{x}) h_k$$

Neural network architecture



Negative log-likelihood loss (probabilistic forecasts)

- In the case of a probabilistic forecast with pdf \hat{f} , we typically measure the prediction error (or loss) by the **negative log-likelihood**

$$\mathcal{L}(y, \hat{f}) = -\ln \hat{f}(y)$$

- We actually never observe a real number y with infinite precision, but an interval $[y]_\epsilon = [y - \epsilon, y + \epsilon]$ centered at y . The probability of that interval is

$$\hat{P}([y]_\epsilon) = \hat{F}(y + \epsilon) - \hat{F}(y - \epsilon) \approx 2\hat{f}(y)\epsilon$$

So, $\mathcal{L}(y, \hat{f}) = -\ln \hat{P}([y]_\epsilon) + \text{cst.}$

- Generalisation to the case of prediction in the form of a belief function?

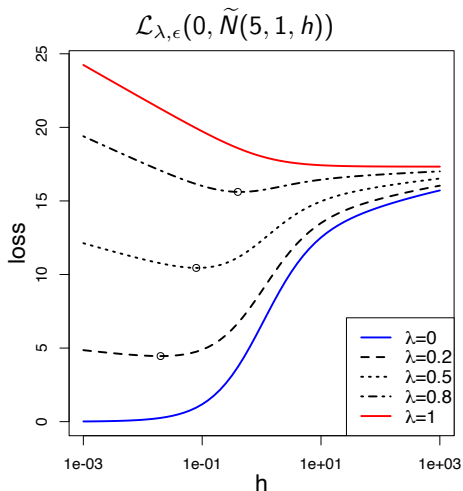
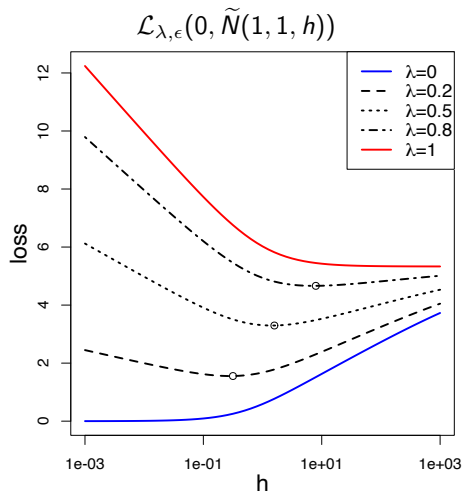
Extension

- $\mathcal{L}_\epsilon(y, \tilde{Y}) = -\ln Bel_{\tilde{Y}}([y]_\epsilon)$ does not work (does not reward imprecision).
- $\mathcal{L}_\epsilon(y, \tilde{Y}) = -\ln Pl_{\tilde{Y}}([y]_\epsilon)$ also does not work (minimized when \tilde{Y} is vacuous).
- Proposal:

$$\mathcal{L}_{\lambda, \epsilon}(y, \tilde{Y}) = -\lambda \ln Bel_{\tilde{Y}}([y]_\epsilon) - (1 - \lambda) \ln Pl_{\tilde{Y}}([y]_\epsilon)$$

with $\lambda \in [0, 1]$ and $\epsilon > 0$.

- Smaller values of λ correspond to more cautious predictions.

Influence of λ 

Training

- The network is trained by minimising the **regularised average loss**

$$C_{\lambda, \epsilon, \xi, \rho}^{(R)}(\Psi) = \underbrace{\frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\lambda, \epsilon}(y_i, \tilde{Y}(\mathbf{x}_i; \Psi))}_{C_{\lambda, \epsilon}(\Psi)} + \underbrace{\frac{\xi}{K} \sum_{k=1}^K h_k}_{R_1(\Psi)} + \underbrace{\frac{\rho}{K} \sum_{k=1}^K \gamma_k^2}_{R_2(\Psi)}$$

where

- $R_1(\Psi)$ has the effect of **reducing the number of prototypes** used for the prediction (setting $h_k = 0$ amounts to discarding prototype k)
 - $R_2(\Psi)$ **shrinks the solution towards a linear model** (setting $\gamma_k = 0$ for all k yields a linear model).
- Heuristics: $\lambda = 0.9$, $\epsilon = 0.01\hat{\sigma}_Y$, ξ and ρ tuned using a validation set or cross-validation.

Calibration

- For any $\alpha \in (0, 1]$, we define an α -level **belief prediction interval (BPI)** as an interval $\mathcal{B}_\alpha(\mathbf{x})$ centered at $\mu(\mathbf{x})$, such that $\text{Bel}_{\tilde{Y}(\mathbf{x})}(\mathcal{B}_\alpha(\mathbf{x})) = \alpha$.
- The predictions are said to be **(weakly) calibrated** if, for all $\alpha \in (0, 1]$, α -level BPIs have a coverage probability at least equal to α , i.e.,

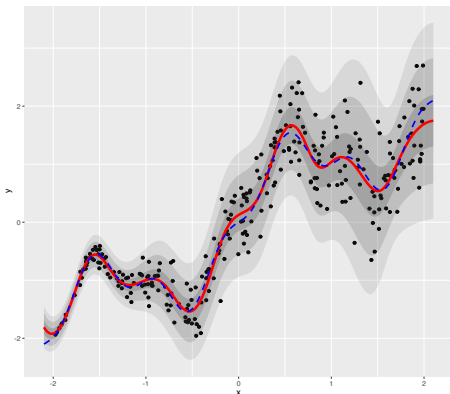
$$\forall \alpha \in (0, 1], \quad P_{\mathbf{X}, Y}(Y \in \mathcal{B}_\alpha(\mathbf{X})) \geq \alpha \quad (1)$$

- As in the probabilistic case, the calibration of evidential predictions can be checked graphically using a **calibration plot** (see infra).
- The precision output $h(\mathbf{x})$ can be multiplied by a constant $c > 0$ to ensure (1) with predictions as precise as possible.

Example

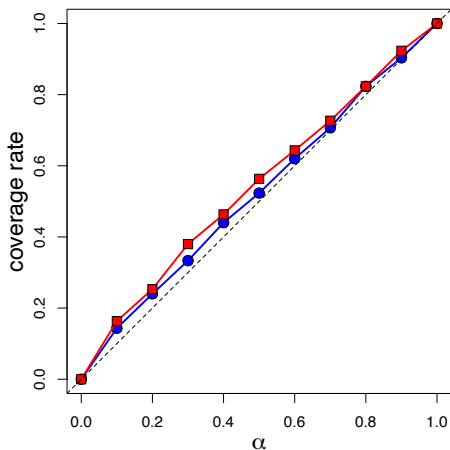
We consider iid data with one-dimensional input $X \sim \text{Unif}(-2, 2)$ and

$$Y = X + (\sin 3X)^3 + \frac{X+2}{4\sqrt{2}}U, \quad U \sim N(0, 1)$$



- Learning and validation sets of size $n = 300$.
- Network with $K = 30$ prototypes initialized by the k-means algorithm.
- ξ and ρ determined by minimizing the validation MSE.
- Shown: true regression function (blue), expected values $\mu(x)$ (red) with BPIs at levels 0.5, 0.9 and 0.99

Calibration curves



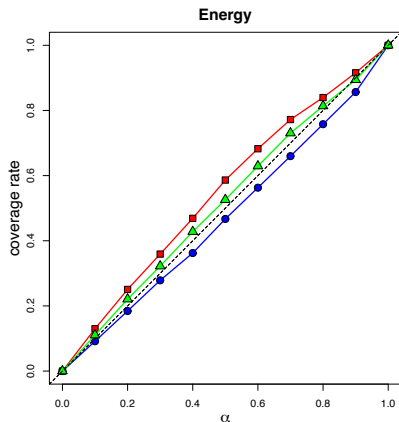
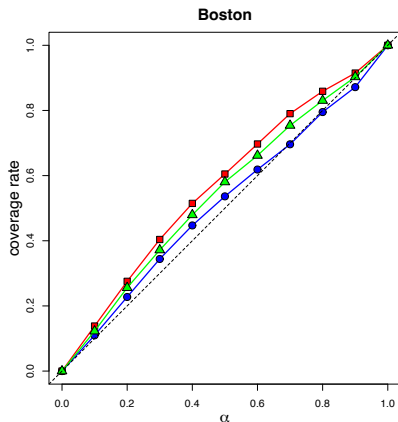
Calibration curves for the probabilistic PIs $\mu(x) \pm u_{(1+\alpha)/2}\sigma(x)$ (in blue) and the BPIs (in red)

Comparison with SOTA methods (RMS & NLL)

	RMS				
	ENNreg	PBP	MC-dropout	Deep ens.	Deep ev. reg.
Boston	2.87 ± 0.14	3.01 ± 0.18	2.97 ± 0.19	3.28 ± 1.00	3.06 ± 0.16
Energy	1.06 ± 0.05	1.80 ± 0.05	1.66 ± 0.04	2.09 ± 0.29	2.06 ± 0.10
Concr.	5.10 ± 0.12	5.67 ± 0.09	5.23 ± 0.12	6.03 ± 0.58	5.85 ± 0.15
Yacht	0.44 ± 0.04	1.02 ± 0.05	1.11 ± 0.09	1.58 ± 0.48	1.57 ± 0.56
Wine	0.63 ± 0.01	0.64 ± 0.01	0.62 ± 0.01	0.64 ± 0.04	0.61 ± 0.02
kin8nm	0.08 ± 0.00	0.10 ± 0.00	0.10 ± 0.00	0.09 ± 0.00	0.09 ± 0.00

	NLL				
	ENNreg	PBP	MC-dropout	Deep ens.	Deep ev. reg.
Boston	2.53 ± 0.07	2.57 ± 0.09	2.46 ± 0.06	2.41 ± 0.25	2.35 ± 0.06
Energy	1.14 ± 0.07	2.04 ± 0.02	1.99 ± 0.02	1.38 ± 0.22	1.39 ± 0.06
Concr.	3.38 ± 0.13	3.16 ± 0.02	3.04 ± 0.02	3.06 ± 0.18	3.01 ± 0.02
Yacht	0.13 ± 0.12	1.63 ± 0.02	1.55 ± 0.03	1.18 ± 0.21	1.03 ± 0.19
Wine	0.94 ± 0.01	0.97 ± 0.01	0.93 ± 0.01	0.94 ± 0.12	0.89 ± 0.05
kin8nm	-1.19 ± 0.00	-0.90 ± 0.01	-0.95 ± 0.01	-1.20 ± 0.02	-1.24 ± 0.01

Calibration plots



Probabilistic predictions (blue), raw evidential predictions (red) and adjusted evidential predictions (green).

Outline

- 1 Distance-based approach
 - Neural network model for classification
 - Neural network model for regression
- 2 Likelihood-based approach
 - Evidential Likelihood-based inference
 - Application to machine learning

Basic ideas

- We now assume the existence of a **parametric statistical model**.
- The **likelihood function** defines a **possibility distribution** in the parameter space, after observing the data.
- Combination of this possibility distribution with random errors yields a **random fuzzy set**, which defines a predictive belief function.
- This approach will be described in general, and then applied to **regression multi-layer neural networks**.

Outline

- 1 Distance-based approach
 - Neural network model for classification
 - Neural network model for regression
- 2 Likelihood-based approach
 - Evidential Likelihood-based inference
 - Application to machine learning

Relative likelihood

- We consider an observed random vector \mathbf{Y} with probability density function (pdf) $f_{\mathbf{Y}|\theta}$, where $\theta \in \Theta$ is the unknown parameter.
- The **likelihood function** after observing $\mathbf{Y} = \mathbf{y}$ is the map

$$L(\cdot; \mathbf{y}) : \Theta \rightarrow \mathbb{R}_+ \\ \theta \mapsto \alpha f_{\mathbf{Y}|\theta}(\mathbf{y})$$

where α is an arbitrary positive constant.

- Assuming that $\sup_{\theta} L(\theta; \mathbf{y}) < +\infty$, we can rescale the likelihood function in $[0, 1]$ and define the **relative likelihood function** as

$$\pi_{\theta|\mathbf{y}} : \theta \mapsto \frac{L(\theta; \mathbf{y})}{\sup_{\theta' \in \Theta} L(\theta'; \mathbf{y})} = \frac{L(\theta; \mathbf{y})}{L(\hat{\theta}; \mathbf{y})}$$

where $\hat{\theta}$ is a maximum likelihood estimate (MLE) of θ .

Possibilistic interpretation

- As remarked by several authors, we can interpret mapping $\pi_{\theta|y} : \Theta \rightarrow [0, 1]$ as a **possibility distribution** over Θ or, equivalently, as the **fuzzy set** of likely values of θ after observing $Y = y$.
- Why does it make sense?
 - Marginalisation**: if $\theta = (\psi, \eta)$, where η is a nuisance parameter, the **profile** relative likelihood is defined as a function of ψ as

$$\pi_{\psi|y}(\psi) = \sup_{\eta} \pi_{\theta|y}(\psi, \eta)$$

which is identical to marginalisation of a possibility distribution.

- Combination**: given two independent observations y and y' ,

$$\pi_{\theta|(y,y')} = \frac{\pi_{\theta|y} \cdot \pi_{\theta|y'}}{\sup_{\theta} \pi_{\theta|y}(\theta) \cdot \pi_{\theta|y'}(\theta)} = \pi_{\theta|y} \cap_p^* \pi_{\theta|y'}$$

where \cap_p^* is the normalised product intersection operation of possibility theory.

Possibilistic interpretation (continued)

- If we accept the possibilistic interpretation of the relative likelihood, the degrees of possibility and necessity of any hypothesis $H \subseteq \Theta$, after observing \mathbf{y} can be computed, respectively, as

$$\Pi_{\theta|\mathbf{y}}(H) = \sup_{\theta \in H} \pi_{\theta|\mathbf{y}}(\theta)$$

and

$$N_{\theta|\mathbf{y}}(H) = 1 - \Pi_{\pi_{\theta|\mathbf{y}}}(H^c)$$

- Under regularity conditions, the large sample distribution of $-2 \ln \Pi_{\theta|\mathbf{y}}(H)$, when H holds, is chi-squared, with df equal to the number r of restrictions imposed by H . Consequently, if H is true,

$$P(\Pi_{\theta|\mathbf{y}}(H) \leq \exp(-\chi_{r;1-\alpha}^2/2)) \approx \alpha$$

i.e., a true hypothesis rarely has a small plausibility (for a large enough sample).

Normal approximation

- Assuming $\ln \pi_{\theta|y}(\theta)$ to be twice differentiable, we can approximate it by its 2nd-order Taylor expansion about $\hat{\theta}$ as

$$\ln \pi_{\theta|y}(\theta) = \ln \pi_{\theta|y}(\hat{\theta}) + (\theta - \hat{\theta})^T \left. \frac{\partial \ln \pi_{\theta|y}(\theta)}{\partial \theta} \right|_{\theta=\hat{\theta}} + \underbrace{\frac{1}{2}(\theta - \hat{\theta})^T \left. \frac{\partial^2 \ln \pi_{\theta|y}(\theta)}{\partial \theta \partial \theta^T} \right|_{\theta=\hat{\theta}}}_{\text{observed information } \mathcal{I}(\hat{\theta})} (\theta - \hat{\theta}) + \dots$$

which gives us the following **normal approximation**:

$$\pi_{\theta|y}(\theta) \approx \exp \left[-\frac{1}{2}(\theta - \hat{\theta})^T \mathcal{I}(\hat{\theta})(\theta - \hat{\theta}) \right]$$

- Thus, $\pi_{\theta|y} \approx \text{GFV}(\hat{\theta}, \mathcal{I}(\hat{\theta}))$.

Prior information

- Assume that we have a Bayesian prior distribution P_0 on Θ .
- Combining P_0 with the $\pi_{\theta|y}$ using the product-intersection rule of GET yields the posterior probability distribution:

$$P_0 \oplus \pi_{\theta|y} = P_{\theta|y}$$

- Evidential likelihood-based inference is, thus, generalised Bayesian inference: it does not require a probabilistic prior, but yields the same result at Bayesian inference is such a prior is provided.
- Prior knowledge may often be more conveniently expressed as a possibility distribution π_0 . In this case, $\pi_0 \oplus \pi_{\theta|y}$ is still a possibility distribution.

Prediction

- We now address the prediction problem: how to **predict/forecast** the value of a r.v. Y_0 with sample space \mathcal{Y} , whose distribution depends on θ ?
- Basic approach:
 - 1 Express Y_0 as a function φ of θ and a pivotal r.v. U with known probability distribution P_U .
 - 2 Combine $\pi_{\theta|y}$, φ and P_U and marginalise on Y_0 . The result is a random fuzzy set.

φ -equation

- We can always write Y_0 as

$$Y_0 = \varphi(\boldsymbol{\theta}, U)$$

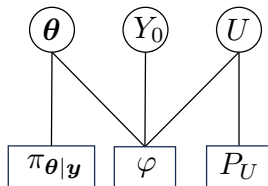
where

- ▶ U is a pivotal r.v. with known distribution and sample space \mathcal{U}
 - ▶ φ is a mapping from $\Theta \times \mathcal{U}$ to \mathcal{Y} .
- Example: if Y_0 is a real r.v. with c.d.f. $F(y; \boldsymbol{\theta})$, we can write

$$Y_0 = F^{-1}(U; \boldsymbol{\theta}) \quad \text{with} \quad U \sim \text{Unif}([0, 1])$$

where F^{-1} is the (generalised) inverse of F .

Predictive random fuzzy set



- Combining $\pi_{\theta|y}$, φ and P_U , and marginalising on Y_0 , we obtain a **random fuzzy set**

$$\begin{aligned} \tilde{Y} : U &\rightarrow \mathcal{F}(\mathcal{Y}) \\ U &\mapsto \pi_{Y_0|y,U} \end{aligned}$$

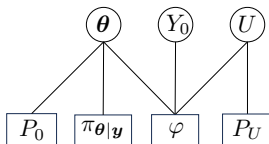
with

$$\pi_{Y_0|y,U}(y) = \sup_{\{\theta \in \Theta : \varphi(\theta, U) = y\}} \pi_{\theta|y}(\theta)$$

- \tilde{Y} induces a **predictive belief function** $Bel_{\tilde{Y}}$.

Special cases

- ① Perfect knowledge of θ : if $\pi_{\theta|y} = \mathbb{1}_{\{\theta_0\}}$, then $Bel_{\tilde{Y}} = P_{Y_0}$.
- ② Bayesian prior:

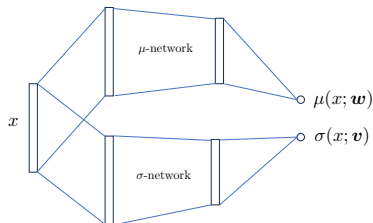


If $\pi_{\theta|y}$ is combined with a Bayesian prior P_0 , then $Bel_{\tilde{Y}} = P_{Y_0|y}$, the Bayesian predictive probability measure.

Outline

- 1 Distance-based approach
 - Neural network model for classification
 - Neural network model for regression
- 2 Likelihood-based approach
 - Evidential Likelihood-based inference
 - Application to machine learning

Model



We consider the following model:

$$Y = \mu(x; \mathbf{w}) + \sigma(x; \mathbf{v})U$$

where

- $U \sim N(0, 1)$ is a random noise with standard normal distribution
- $\mu(x; \mathbf{w})$, $\sigma(x; \mathbf{v})$ are the outputs of two neural networks with distinct weight vectors \mathbf{w} and \mathbf{v} , referred to, respectively, as the μ and σ -networks

Log-likelihood

- Let $\boldsymbol{\theta} = (\mathbf{w}^T, \mathbf{v})^T$ denote the vector of all parameters in the model. Given i.i.d. data $\{(x_i, y_i)\}_{i=1}^n$, the log-likelihood of $\boldsymbol{\theta}$ is

$$\ell(\boldsymbol{\theta}) = - \sum_{i=1}^n \ln \sigma(x_i, \mathbf{v}) - \frac{1}{2} \ln(2\pi) - \frac{1}{2} \sum_{i=1}^n \left(\frac{y_i - \mu(x_i; \mathbf{w})}{\sigma(x_i, \mathbf{v})} \right)^2$$

- We fit the model by maximising the **regularised log-likelihood**

$$\ell_{\boldsymbol{\lambda}, \boldsymbol{\xi}}^{(R)}(\boldsymbol{\theta}) = \ell(\boldsymbol{\theta}) - \sum_{j=1}^{N_w} \lambda_j w_j^2 - \sum_{j=1}^{N_v} \xi_j v_j^2$$

where $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_{N_w})$ and $\boldsymbol{\xi} = (\xi_1, \dots, \xi_{N_v})$ are vectors of regularisation coefficients.

Relative log-likelihood

- Let $\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{w}}^T, \hat{\boldsymbol{v}}^T)^T$ be a **global maximiser** of $\ell_{\boldsymbol{\lambda}, \boldsymbol{\xi}}^{(R)}(\boldsymbol{\theta})$. We define the possibility distribution of $\boldsymbol{\theta}$ as

$$\pi_{\boldsymbol{\theta}|\boldsymbol{y}}(\boldsymbol{\theta}) = \exp \left[\ell_{\boldsymbol{\lambda}, \boldsymbol{\xi}}^{(R)}(\boldsymbol{\theta}) - \ell_{\boldsymbol{\lambda}, \boldsymbol{\xi}}^{(R)}(\hat{\boldsymbol{\theta}}) \right]$$

- We can write

$$\pi_{\boldsymbol{\theta}|\boldsymbol{y}}(\boldsymbol{\theta}) \propto \underbrace{\exp(\ell(\boldsymbol{\theta}) - \ell(\hat{\boldsymbol{\theta}}))}_{\text{relative likelihood}} \times \underbrace{\exp\left(-\boldsymbol{\theta}^T \text{diag}(\boldsymbol{\lambda}, \boldsymbol{\xi})\boldsymbol{\theta}\right)}_{\text{prior } \pi_0(\boldsymbol{\theta})}$$

where $\pi_0 \sim \text{GFV}(\mathbf{0}, 2 \text{diag}(\boldsymbol{\lambda}, \boldsymbol{\xi}))$.

Normal approximation

- Using the normal approximation of the relative likelihood, we can write

$$\pi_{\theta|y} \sim \text{GFV} \left(\hat{\theta}, \mathcal{I}_{\lambda, \xi}(\hat{\theta}) \right)$$

with precision matrix

$$\mathcal{I}_{\lambda, \xi}(\hat{\theta}) = -\frac{\partial^2 \ell_{\lambda, \xi}^{(R)}}{\partial \theta \partial \theta^T}(\hat{\theta})$$

- This matrix can be computed by backpropagation or by numerical differentiation (slower, but this calculation needs to be performed only once).

φ -equation and linearisation

- We now wish to predict the **response** Y_0 for a new input vector $x = x_0$; it can be written as

$$Y_0 = \mu(x_0, \mathbf{w}) + \sigma(x_0, \mathbf{v})U = \varphi(x_0, \boldsymbol{\theta}, U)$$

with $U \sim N(0, 1)$.

- To approximate the condition possibility distribution $\pi_{Y_0|y,u}$ of Y_0 given $U = u$, we need to **linearise** both $\mu(x_0; \mathbf{w})$ and $\sigma(x_0, \mathbf{v})$ around, respectively, $\hat{\mathbf{w}}$ and $\hat{\mathbf{v}}$. We write

$$\mu(x_0; \mathbf{w}) \approx \mu(x_0; \hat{\mathbf{w}}) + \nabla \mu_0^T (\mathbf{w} - \hat{\mathbf{w}})$$

$$\sigma(x_0; \mathbf{v}) \approx \sigma(x_0; \hat{\mathbf{v}}) + \nabla \sigma_0^T (\mathbf{v} - \hat{\mathbf{v}})$$

with

$$\nabla \mu_0 = \left. \frac{\partial \mu(x_0; \mathbf{w})}{\partial \mathbf{w}} \right|_{\mathbf{w}=\hat{\mathbf{w}}} \quad \text{and} \quad \nabla \sigma_0 = \left. \frac{\partial \sigma(x_0; \mathbf{v})}{\partial \mathbf{v}} \right|_{\mathbf{v}=\hat{\mathbf{v}}}$$

Predictive RFS

- We thus have

$$Y_0 \approx \begin{pmatrix} \nabla \mu_0 \\ U \nabla \sigma_0 \end{pmatrix}^T \boldsymbol{\theta} + \mu(x_0, \hat{\mathbf{w}}) - \nabla \mu_0^T \hat{\mathbf{w}} + (\sigma(x_0, \hat{\mathbf{v}}) - \nabla \sigma_0^T \hat{\mathbf{v}}) U$$

- Assuming matrix $\mathcal{I}_{\lambda, \xi}(\hat{\boldsymbol{\theta}})$ to be positive definite, and using a result about affine transformations of GRVs, it can be shown that, approximately,

$$\pi_{Y_0|y, u} \sim \text{GFN} \left(\mu(x_0, \hat{\mathbf{w}}) + \sigma(x_0, \hat{\mathbf{v}}) u, h(x_0, \hat{\boldsymbol{\theta}}, u) \right)$$

with

$$h(x_0, \hat{\boldsymbol{\theta}}, u) = \left[\begin{pmatrix} \nabla \mu_0^T & u \nabla \sigma_0 \end{pmatrix} \mathcal{I}_{\lambda, \xi}(\hat{\boldsymbol{\theta}})^{-1} \begin{pmatrix} \nabla \mu_0 \\ u \nabla \sigma_0 \end{pmatrix} \right]^{-1} = \frac{1}{\alpha + \gamma u + \delta u^2}$$

- Uncertainty about Y_0 can, thus, be described by the **predictive RFS**

$$\tilde{Y}(x_0) : U \mapsto \pi_{Y_0|y, U} \sim \text{GFN} \left(\mu(x_0; \hat{\mathbf{w}}) + U \sigma(x_0; \hat{\mathbf{v}}), h(x_0, \hat{\boldsymbol{\theta}}, U) \right)$$

GRFN approximation

- From Jensen's inequality,

$$\mathbb{E} \left[\frac{1}{\alpha + \gamma U + \delta U^2} \right] \geq \frac{1}{\alpha + \gamma \mathbb{E}[U] + \delta \mathbb{E}[U^2]} = \frac{1}{\alpha + \delta}$$

- Replacing $h(x_0, \hat{\theta}, U)$ by the lower bound of its expectation

$$\tilde{h}(x_0, \hat{\theta}) = \frac{1}{\alpha + \delta}$$

we obtain the conservative approximation:

$$\tilde{Y}(x_0) \sim \tilde{N} \left(\mu(x_0; \hat{\mathbf{w}}), \sigma^2(x_0; \hat{\mathbf{v}}), \tilde{h}(x_0, \hat{\theta}) \right)$$

Experimental results

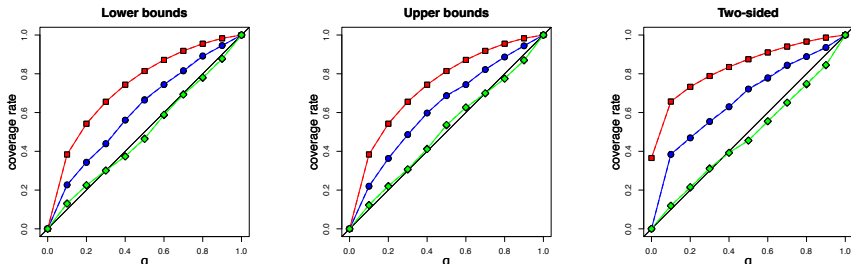
- Real datasets:

	n	p	response
Boston	506	13	medv
Energy	768	8	Y2
Concrete	1030	8	strength
Yacht	308	5	Y
Airfoil	1503	5	Y

- Neural networks with ELU activation functions, one layer of 5 units for the σ -network, one layer of 30 or 50 units for the μ -network.

Calibration curves

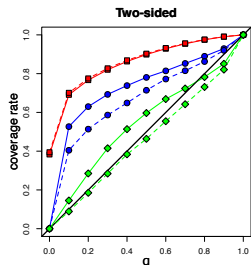
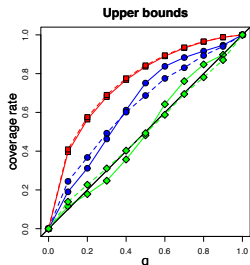
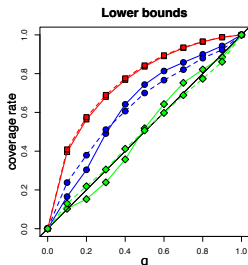
Boston dataset



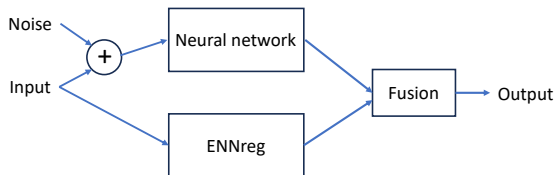
Calibration curves for belief lower bounds (left), belief upper bounds (middle) and two-sided belief intervals (right) computed using the validation data. The lower green curves, middle blue curve and upper red curves correspond, respectively, to the coverage rates of probabilistic predictive intervals, coverage rates of belief intervals, and average plausibilities of belief intervals

Calibration curves

Energy dataset

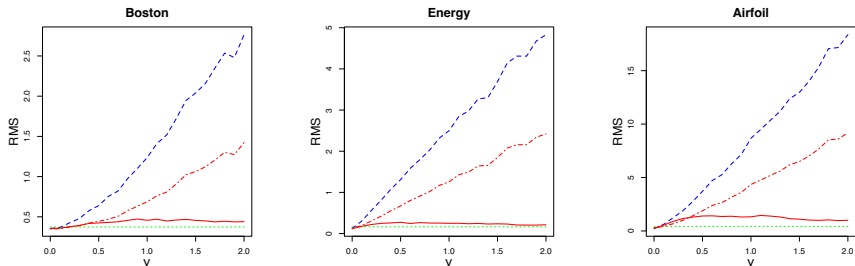


Experiment 1



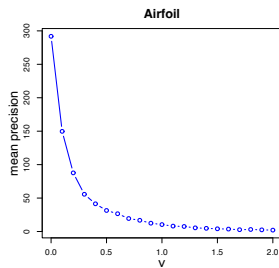
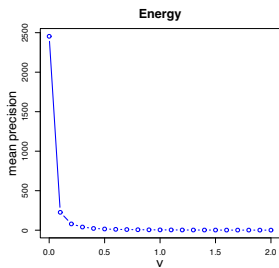
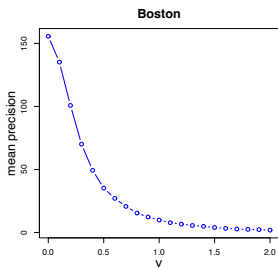
- To each test input of the neural network, we added a Gaussian noise with mean 0 and standard deviation ν .
- The outputs of the two models were combined by
 - ▶ averaging (using only the point predictions)
 - ▶ the product-intersection rule.
- The process was repeated ten times for each noise level, with ten 80%-20% training/test partitions.

Results



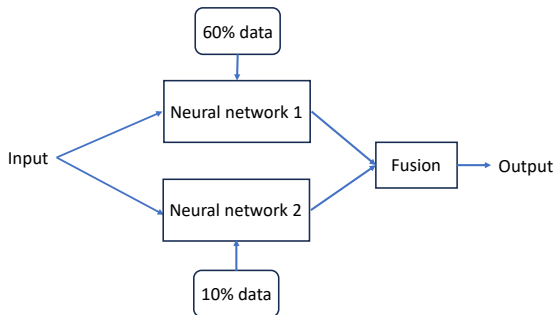
Root-mean-squared (RMS) error as a function of noise standard deviation ν for the following methods: (1) neural network with corrupted inputs (blues broken line), (2) ENNreg with uncorrupted inputs (horizontal green dotted line), average prediction (red dash-dotted line), (3) combination by the product-intersection rule (solid red line).

Precisions



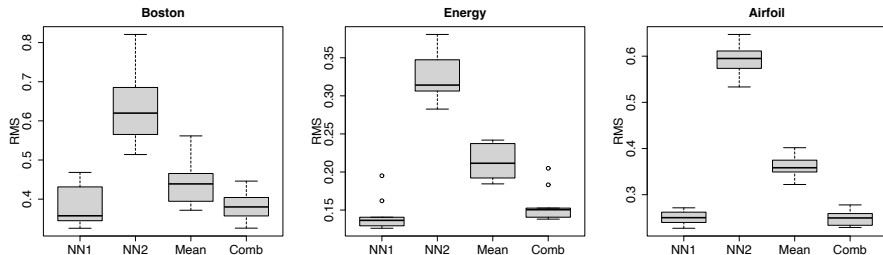
Mean precision as a function of noise standard deviation ν for the five datasets

Experiment 2



- The reliability of predictions also depends on the size of the training set.
- To simulate another fusion scenario involving the combination of predictions with different reliabilities, we combined the outputs of two neural networks, trained with 60% and 10% of the data.
- The process was repeated ten times with different random training/validation/test partitions.

Results



Box plots of root-mean-squared errors for the five datasets and the following methods: neural network trained on 60% of the data (NN1), neural network trained on 10% of the data (NN2), average prediction (Mean) and combination by the product-intersection rule (Comb).

Summary

- **Evidential ML** = ML with uncertainty quantified by belief functions.
- Two main approaches:
 - ① Distance-based (nonparametric)
 - ② Likelihood-based (parametric)
- The distance-based approach, initially proposed for classification, can be extended to regression using GRFNs.
- The likelihood-based approach yields predictive RFSs. It can be applied to regression or classification.
- These methods are based on **generalised evidence theory**, a sound and intuitively appealing theoretical framework.
- They provide **cautious predictions**, especially useful in situations of high uncertainty:
 - ▶ Small datasets
 - ▶ Low-quality data
 - ▶ Several sources providing partial information

References on ENN

<https://www.hds.utc.fr/~tdenoeux>



T. Denœux.

A neural network classifier based on Dempster-Shafer theory

IEEE Transactions on Systems, Man and Cybernetics A, 30(2):131-150, 2000



Z. Tong, Ph. Xu and T. Denœux.

An evidential classifier based on Dempster-Shafer theory and deep learning

Neurocomputing 450:275–293, 2021



L. Huang, S. Ruan, P. Decazes and T. Denœux

Lymphoma segmentation from 3D PET-CT images using a deep evidential network

International Journal of Approximate Reasoning 149:39-60, 2022



L. Huang, S. Ruan, P. Decazes and T. Denœux

Deep evidential fusion with uncertainty quantification and reliability learning for multimodal medical image segmentation

Information Fusion, 113:102648, 2025

References on ENNreg

<https://www.hds.utc.fr/~tdenoeux>



T. Denœux

Quantifying Prediction Uncertainty in Regression using Random Fuzzy Sets:
the ENNreg model

IEEE Transactions on Fuzzy Systems, 31(10):3690–3699, 2023






L. Huang, Y.Xing, S. Mishra, T. Denœux and M. Feng.

Evidential time-to-event prediction with calibrated uncertainty quantification

International Journal of Approximate Reasoning, 181:109403, 2025

References on the likelihood-based approach

<https://www.hds.utc.fr/~tdenoeux>

-  O. Kanjanatarakul, S. Sriboonchitta and T. Denœux.
Forecasting using belief functions: an application to marketing econometrics.
International Journal of Approximate Reasoning, 55(5):1113–1128, 2014.
-  O. Kanjanatarakul, T. Denœux and S. Sriboonchitta.
Prediction of future observations using belief functions: a likelihood-based approach
International Journal of Approximate Reasoning 72:71–94, 2016
-  T. Denœux
Uncertainty Quantification in Logistic Regression using Random Fuzzy Sets and Belief Functions
International Journal of Approximate Reasoning, Vol. 168, 109159, 2024
-  T. Denœux
Uncertainty Quantification in Regression Neural Networks using Evidential Likelihood-based Inference
International Journal of Approximate Reasoning, Vol. 182, 109423, 2025