

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

Thierry Denoeux <sup>1</sup>

<sup>1</sup>Université de technologie de Compiègne

October 30, 2023

## Abstract

We introduce a neural network model for regression in which prediction uncertainty is quantified by Gaussian random fuzzy numbers (GRFNs), a newly introduced family of random fuzzy subsets of the real line that generalizes both Gaussian random variables and Gaussian possibility distributions. The output GRFN is constructed by combining GRFNs induced by prototypes using a combination operator that generalizes Dempster's rule of Evidence Theory. The three output units indicate the most plausible value of the response variable, variability around this value, and epistemic uncertainty. The network is trained by minimizing a loss function that generalizes the negative log-likelihood. Comparative experiments show that this method is competitive, both in terms of prediction accuracy and calibration error, with state-of-the-art techniques such as random forests or deep learning with Monte Carlo dropout. In addition, the model outputs a predictive belief function that can be shown to be calibrated, in the sense that it allows us to compute conservative prediction intervals with specified belief degree.

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

Thierry Denœux

**Abstract**—We introduce a neural network model for regression in which prediction uncertainty is quantified by Gaussian random fuzzy numbers (GRFNs), a newly introduced family of random fuzzy subsets of the real line that generalizes both Gaussian random variables and Gaussian possibility distributions. The output GRFN is constructed by combining GRFNs induced by prototypes using a combination operator that generalizes Dempster’s rule of Evidence Theory. The three output units indicate the most plausible value of the response variable, variability around this value, and epistemic uncertainty. The network is trained by minimizing a loss function that generalizes the negative log-likelihood. Comparative experiments show that this method is competitive, both in terms of prediction accuracy and calibration error, with state-of-the-art techniques such as random forests or deep learning with Monte Carlo dropout. In addition, the model outputs a predictive belief function that can be shown to be calibrated, in the sense that it allows us to compute conservative prediction intervals with specified belief degree.

**Index Terms**—Evidence theory, Dempster-Shafer theory, belief functions, machine learning, neural networks.

## I. INTRODUCTION

The Dempster-Shafer (DS) theory of belief functions, also known as Evidence Theory, is a mathematical formalism for reasoning with uncertainty [1], [2]. This formalism relies on two main components: (1) the representation of elementary pieces of evidence using belief functions, and (2) their combination by a conjunctive operator, referred to as Dempster’s rule of combination. The greater flexibility of DS theory makes it suitable to reason and make decisions based on weak information, which would not be adequately represented by probability measures. In particular, it is possible, using belief functions, to distinguish between conflicting evidence equally supporting different hypotheses on the one hand, and sheer ignorance resulting from total lack of evidence on the other hand [3].

In machine learning, DS theory provides a powerful framework for expressing uncertainty about the output of a learning algorithm. Most applications so far have concerned clustering [4], [5] and classification [6]–[8]. In particular, in supervised classification, several algorithms have been proposed to learn *evidential classifiers*, defined as classifiers that represent class prediction uncertainty by belief functions. One of the first such classifier is the evidential  $K$ -nearest neighbor (EKNN) rule [6], which constructs a predictive belief function by combining the evidence of the  $K$  nearest neighbors of the feature vector under consideration using Dempster’s rule. A similar idea is implemented in the evidential neural network

(ENN) model introduced in [9], in which the learning vectors are summarized by prototypes, whose location in feature space is optimized together with other network parameters. Recently, this idea has been applied to deep networks [10] [11] by adding a “DS layer” to a deep architecture.

The reason why applications of DS theory to machine learning have been confined to clustering and classification is that these learning tasks only require the definition of belief functions over a *finite* set of clusters or classes, a simple mathematical framework for which many tools and concepts of DS theory have been developed [2]. A belief function on a finite domain (or *frame of discernment*) can be conveniently represented by a mass function assigning probability masses to some subsets called focal sets. Such mass functions are easily combined by Dempster’s rule. In contrast, applying DS theory to regression is more challenging, as the domain of the response variable is then the real line or a real interval. Whereas some mathematical models for generating belief functions on the real line, such as random intervals have been well studied [12], [13], combining such belief functions by Dempster’s rule is often computationally difficult and requires Monte Carlo simulation [14].

A way to circumvent this difficulty is to discretize the response variable. This is the strategy followed in [15], in which a regression neural network directly extending the ENN model is proposed. The output of this model is a mass function focussed on disjoint intervals and the whole domain of the response variable. However, discretizing a variable inevitably results in a loss of accuracy. Another approach, introduced in [16], [17], is to directly extend the EKNN rule by defining, for each neighbor, a simple mass functions with two focal sets: a singleton composed of a real number and the frame of discernment. This simple method, called EVREG, was shown in [17] to have good performances as compared to nearest neighbor regression and other nonparametric methods, and to allow efficient handling of uncertain response data (such as provided by an unreliable sensor). However, as a nonparametric method, the performance of EVREG decreases with the input dimension, and it cannot compete with state-of-the-art regression methods.

In this paper, we propose a new regression neural network model that quantifies prediction uncertainty using belief functions<sup>1</sup>. This model, called ENNreg, is based on the formalism of *Gaussian random fuzzy numbers (GRFNs)* recently introduced in [19]. A GRFN is a random fuzzy subset of

T. Denœux is with Université de technologie de Compiègne, CNRS, Heudiasyc, France, and with Institut universitaire de France, Paris, France.

<sup>1</sup>This paper is an extended version of the short paper [18] presented at the 7th International Conference on Belief Functions, Paris, 26-28 October 2022.

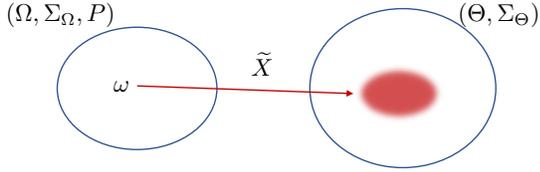


Fig. 1: Definition of a random fuzzy set.

the real line, which can be described as a Gaussian possibility distribution whose mode is a Gaussian random variable or, equivalently, as a Gaussian random variable with fuzzy mean but precise variance. GRFNs induce belief functions on the real line and can be combined using a generalization of Dempster’s rule. In ENNreg, a GRFN is associated to each prototype based on the distance to the input vector, and the GRFNs are combined using an operator generalizing Dempster’s rule. The network output is a GRFN quantifying both random and epistemic prediction uncertainty. The method is shown to be competitive with state-of-the-art models in a variety of regression tasks, while quantifying prediction uncertainty using a GRFN and its associated belief function. The output belief function is shown to be *calibrated*, in the sense that it makes it possible to compute prediction intervals with coverage probability at least as large as their nominal level.

The rest of this paper is organized as follows. The necessary background on random fuzzy sets and GRFNs is first recalled in Section II. The desired properties of predictive belief functions in regression tasks are then studied in Section III, and the proposed ENNreg model is introduced in Section IV. Experimental results are reported in Section V, and Section VI concludes the paper.

## II. RANDOM FUZZY SETS

The theory of epistemic random fuzzy sets (RFSs) was introduced in [20] and [19] as a general framework encompassing both DS theory and possibility theory. We first recall some important definitions related to random fuzzy sets in Section II-A. Gaussian random fuzzy numbers, a parametric family of RFSs on the real line are then described in Section II-B.

### A. General definitions

The notion of random fuzzy set generalizes that of random set [21] by mapping elements of a probability space to *fuzzy subsets* of the domain  $\Theta$  of the variable of interest. This mathematical framework has been used to model a random experiment whose outcomes are fuzzy sets [22]–[24]. Here, as in [20] and [19], we use it as a model of unreliable and fuzzy evidence, directly generalizing the DS framework.

Formally, let  $(\Omega, \Sigma_\Omega, P)$  be a probability space,  $(\Theta, \Sigma_\Theta)$  a measurable space, and  $\tilde{X}$  a mapping from  $\Omega$  to the set  $[0, 1]^\Theta$  of *fuzzy subsets* of  $\Theta$  (see Figure 1). Under measurability conditions [25], the tuple  $(\Omega, \Sigma_\Omega, P, \Theta, \Sigma_\Theta, \tilde{X})$  is said to be a *random fuzzy set*.

Under the epistemic interpretation, we see  $\Omega$  as a *set of interpretations* of a piece of evidence about a variable  $\theta$

taking values in  $\Theta$ . If interpretation  $\omega \in \Omega$  holds, we know that “ $\theta$  is  $\tilde{X}(\omega)$ ”, i.e.,  $\theta$  is constrained by the possibility distribution defined by  $\tilde{X}(\omega)$ . The resulting random fuzzy set thus encodes a state of knowledge about variable  $\theta$ . If all images  $\tilde{X}(\omega)$  are crisp, then  $\tilde{X}$  defines an ordinary random set. If  $\tilde{X}$  is a constant mapping, it is equivalent to specifying a unique fuzzy subset of  $\Theta$ , which defines a possibility distribution [26].

*Belief and plausibility functions:* In the following, we assume random fuzzy set  $\tilde{X}$  to be *normalized*, i.e., to verify the following conditions:

- 1) For any  $\omega \in \Omega$ ,  $\tilde{X}(\omega)$  is either the empty set, or a normal fuzzy set, i.e., the height of  $\tilde{X}(\omega)$ , defined as  $\text{hgt}(\tilde{X}(\omega)) = \sup_{\theta \in \Theta} \tilde{X}(\omega)(\theta)$ , is either 0 or 1;
- 2) The set of elements  $\omega \in \Omega$  whose image is empty has zero probability, i.e.,  $P(\{\omega \in \Omega : \tilde{X}(\omega) = \emptyset\}) = 0$ .

Conditionally on  $\omega \in \Omega$  being the true interpretation, we can define possibility and necessity measures [26] on  $\Theta$ , respectively, as

$$\Pi_{\tilde{X}(\omega)}(B) = \sup_{\theta \in B} \tilde{X}(\omega)(\theta), \quad (1)$$

and

$$N_{\tilde{X}(\omega)}(B) = \begin{cases} 1 - \Pi_{\tilde{X}(\omega)}(B^c) & \text{if } \tilde{X}(\omega) \neq \emptyset \\ 0 & \text{otherwise,} \end{cases}$$

for all  $B \subseteq \Theta$ . Let  $Bel_{\tilde{X}}$  and  $Pl_{\tilde{X}}$  denote the mappings from  $\Sigma_\Theta$  to  $[0, 1]$  that associate to each subset  $B \in \Sigma_\Theta$ , respectively, its expected necessity and its expected possibility, i.e.,

$$Bel_{\tilde{X}}(B) = \int_{\Omega} N_{\tilde{X}(\omega)}(B) dP(\omega) \quad (2)$$

and

$$Pl_{\tilde{X}}(B) = \int_{\Omega} \Pi_{\tilde{X}(\omega)}(B) dP(\omega) = 1 - Bel_{\tilde{X}}(B^c). \quad (3)$$

It can be shown that  $Bel_{\tilde{X}}$  and  $Pl_{\tilde{X}}$  are, respectively, belief and plausibility functions [25].

*Combination:* Consider two epistemic random fuzzy sets  $(\Omega_i, \Sigma_i, P_i, \Theta, \Sigma_\Theta, \tilde{X}_i)$ ;  $i = 1, 2$ , encoding independent pieces of evidence. If interpretations  $\omega_1 \in \Omega_1$  and  $\omega_2 \in \Omega_2$  both hold, we know that “ $\theta$  is  $\tilde{X}_1(\omega_1)$ ” and “ $\theta$  is  $\tilde{X}_2(\omega_2)$ ”. It is then natural to combine the fuzzy sets  $\tilde{X}_1(\omega_1)$  and  $\tilde{X}_2(\omega_2)$  by an intersection operator. As discussed in [20], the normalized product intersection operator  $\odot$ , defined for any two fuzzy subsets  $\tilde{F}$  and  $\tilde{G}$  of  $\Theta$  as

$$(\tilde{F} \odot \tilde{G})(\theta) = \begin{cases} \frac{\tilde{F}(\theta)\tilde{G}(\theta)}{\text{hgt}(\tilde{F} \cdot \tilde{G})} & \text{if } \text{hgt}(\tilde{F} \cdot \tilde{G}) > 0 \\ 0 & \text{otherwise,} \end{cases} \quad (4)$$

is suitable for combining fuzzy information from independent sources and it is associative [27]. We thus consider the mapping  $\tilde{X}_\odot$  from  $\Omega_1 \times \Omega_2$  to  $[0, 1]^\Theta$  defined by  $\tilde{X}_\odot(\omega_1, \omega_2) = \tilde{X}_1(\omega_1) \odot \tilde{X}_2(\omega_2)$ . To insure the normality condition, we need to condition the product measure  $P_1 \times P_2$  on the set  $\Theta_{12}^*$  of (partially) consistent interpretations, defined as

$$\Theta_{12}^* = \{(\omega_1, \omega_2) \in \Omega_1 \times \Omega_2 : \tilde{X}_\odot(\omega_1, \omega_2) \neq \emptyset\}.$$

We then obtain the combined random fuzzy set

$$(\Omega_1 \times \Omega_2, \Sigma_1 \otimes \Sigma_2, P_{12}, \Theta, \Sigma_\Theta, \tilde{X}_\Theta),$$

where  $P_{12}$  is the conditioned product measure. This operation will be referred to as the *generalized product-intersection rule with hard normalization*<sup>2</sup>. The corresponding operator, denoted by  $\boxplus$ , is commutative and associative, and it is a proper generalization of Dempster's rule (it boils down to Dempster's rule when the combined random fuzzy sets are crisp).

### B. Gaussian Random Fuzzy Numbers

To define a practical model of random fuzzy subset of the real line, we start with the definition of a *Gaussian Fuzzy Number* (GFN). A GFN is a fuzzy subset of  $\mathbb{R}$  with membership function

$$\varphi(x; m, h) = \exp\left(-\frac{h}{2}(x - m)^2\right),$$

where  $m \in \mathbb{R}$  is the mode and  $h \in [0, +\infty]$  is the precision. Such a fuzzy number will be denoted by  $\text{GFN}(m, h)$ . An interesting property of this family of fuzzy numbers is that it is closed with respect to the normalized product intersection: given two fuzzy numbers  $\text{GFN}(m_1, h_1)$  and  $\text{GFN}(m_2, h_2)$ , their combination by the  $\odot$  operator (4) yields the fuzzy number  $\text{GFN}(m_{12}, h_{12})$  with

$$m_{12} = \frac{h_1 m_1 + h_2 m_2}{h_1 + h_2} \quad \text{and} \quad h_{12} = h_1 + h_2. \quad (5)$$

A *Gaussian random fuzzy number* (GRFN) can now be defined as a GFN whose mode is a Gaussian random variable. More formally, let  $(\Omega, \Sigma_\Omega, P)$  be a probability space and let  $M : \Omega \rightarrow \mathbb{R}$  be a Gaussian random variable (GRV) with mean  $\mu$  and variance  $\sigma^2$ . The random fuzzy set  $\tilde{X} : \Omega \rightarrow [0, 1]^{\mathbb{R}}$  defined as

$$\tilde{X}(\omega) = \text{GFN}(M(\omega), h)$$

is called a GRFN with mean  $\mu$ , variance  $\sigma^2$  and precision  $h$ , which we write  $\tilde{X} \sim \tilde{N}(\mu, \sigma^2, h)$ . A GRFN is, thus, defined by a location parameter  $\mu$ , and two parameters  $h$  and  $\sigma^2$  corresponding, respectively, to possibilistic and probabilistic uncertainty.

A GRFN can be seen either as a generalized GRV with fuzzy mean, or as a generalized GFN with random mode. In particular, a GRFN  $\tilde{X}$  with infinite precision  $h = +\infty$  is equivalent to a GRV with mean  $\mu$  and variance  $\sigma^2$ , which we can write:  $\tilde{N}(\mu, \sigma^2, +\infty) = N(\mu, \sigma^2)$ . If  $\sigma^2 = 0$ ,  $M$  is a constant random variable taking value  $\mu$ , and  $\tilde{X}$  is a possibilistic variable with possibility distribution  $\text{GFN}(\mu, h)$ . Another case of interest is that where  $h = 0$ , in which case  $\tilde{X}(\omega)(x) = 1$  for all  $\omega \in \Omega$  and  $x \in \mathbb{R}$ ; the belief function induced by  $\tilde{X}$  is then vacuous.

As an illustration, Figure 2 shows ten realizations of two GRFNs  $\tilde{X}_1 \sim \tilde{N}(0, (0.2)^2, 0.2)$  and  $\tilde{X}_2 \sim \tilde{N}(5, 1, 5)$ . GRFN  $\tilde{X}_1$  (left) has high imprecision and low variability, while  $\tilde{X}_2$  (right) is more precise but has higher variability.

<sup>2</sup>In [19], we proposed a different operation based on a ‘‘soft normalization’’ process consisting in conditioning the product measure by a fuzzy subset consistent interpretations. Here, hard normalization is preferred as it results in much simpler calculations.

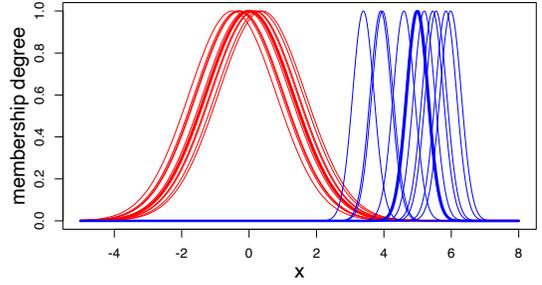


Fig. 2: Ten realizations of GRFNs  $\tilde{X}_1 \sim \tilde{N}(0, (0.2)^2, 0.2)$  (left, in red) and  $\tilde{X}_2 \sim \tilde{N}(5, 1, 5)$  (right, in blue). For each GRFN  $\tilde{X}_k$ , the possibility distribution  $\text{GFN}(\mu_k, h_k)$  is also plotted as a thick line.

As shown in [19], the contour function  $pl_{\tilde{X}}$  of  $\tilde{X}$ , defined as  $pl_{\tilde{X}}(x) = Pl_{\tilde{X}}(\{x\})$  for all  $x \in \mathbb{R}$ , has the following expression:

$$pl_{\tilde{X}}(x) = \frac{1}{\sqrt{1 + h\sigma^2}} \exp\left(-\frac{h(x - \mu)^2}{2(1 + h\sigma^2)}\right), \quad (6)$$

and the belief and plausibility of any real interval  $[x, y]$  are given by the following formulas:

$$\begin{aligned} Bel_{\tilde{X}}([x, y]) &= \Phi\left(\frac{y - \mu}{\sigma}\right) - \Phi\left(\frac{x - \mu}{\sigma}\right) - \\ &pl_{\tilde{X}}(x) \left[ \Phi\left(\frac{(x + y)/2 - \mu}{\sigma\sqrt{h\sigma^2 + 1}}\right) - \Phi\left(\frac{x - \mu}{\sigma\sqrt{h\sigma^2 + 1}}\right) \right] - \\ &pl_{\tilde{X}}(y) \left[ \Phi\left(\frac{y - \mu}{\sigma\sqrt{h\sigma^2 + 1}}\right) - \Phi\left(\frac{(x + y)/2 - \mu}{\sigma\sqrt{h\sigma^2 + 1}}\right) \right], \quad (7a) \end{aligned}$$

and

$$\begin{aligned} Pl_{\tilde{X}}([x, y]) &= \Phi\left(\frac{y - \mu}{\sigma}\right) - \Phi\left(\frac{x - \mu}{\sigma}\right) + \\ &pl_{\tilde{X}}(x) \Phi\left(\frac{x - \mu}{\sigma\sqrt{h\sigma^2 + 1}}\right) + \\ &pl_{\tilde{X}}(y) \left[ 1 - \Phi\left(\frac{y - \mu}{\sigma\sqrt{h\sigma^2 + 1}}\right) \right], \quad (7b) \end{aligned}$$

where  $\Phi$  is the standard normal cumulative distribution function (cdf).

The usefulness of GRFNs as a model of uncertain information about a real quantity arises from the fact that GRFNs can easily be combined by the generalized product-intersection rule  $\boxplus$  introduced in II-A. Given two GRFNs  $\tilde{X}_1 \sim \tilde{N}(\mu_1, \sigma_1^2, h_1)$  and  $\tilde{X}_2 \sim \tilde{N}(\mu_2, \sigma_2^2, h_2)$ , we have, as a direction consequence of (5),  $\tilde{X}_1 \boxplus \tilde{X}_2 \sim \tilde{N}(\mu_{12}, \sigma_{12}^2, h_{12})$ , with

$$\mu_{12} = \frac{h_1 \mu_1 + h_2 \mu_2}{h_1 + h_2}, \quad \sigma_{12}^2 = \frac{h_1^2 \sigma_1^2 + h_2^2 \sigma_2^2}{(h_1 + h_2)^2}, \quad (8a)$$

$$\text{and} \quad h_{12} = h_1 + h_2. \quad (8b)$$

### III. PREDICTIVE RANDOM FUZZY SETS

In this section and in the rest of this paper, we consider a regression problem, in which the task is to predict a continuous random response variable  $Y$  from a random

vector  $\mathbf{X} = (X_1, \dots, X_p)$  of  $p$  input variables. In Section IV we will propose a neural network model that, given an observed input vector  $\mathbf{X} = \mathbf{x}$ , computes a GRFN  $\tilde{Y}(\mathbf{x})$  with associated belief function  $Bel_{\tilde{Y}(\mathbf{x})}$  that ‘‘approximates’’  $Y$  in some way. The network will be trained using a training set  $\mathcal{T} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  assumed to be an independent and identically distributed (iid) sample from  $(\mathbf{X}, Y)$ . In this section, we define desirable properties for a predictive RFS  $\tilde{Y}(\mathbf{x})$ . A loss function measuring predictive accuracy is first defined in Section III-A. A notion of calibration is then introduced in Section III-B.

### A. Loss function

When predicting a random variable  $Y$  by a probability measure  $\hat{P}$  with density function  $\hat{f}$ , we typically measure the prediction error (or loss) by the negative log-likelihood

$$\mathcal{L}(y, \hat{f}) = -\ln \hat{f}(y), \quad (9)$$

which can be interpreted as follows. The continuous variable  $Y$  is always observed with finite precision, so we actually observe an interval  $[y]_\epsilon = [y - \epsilon, y + \epsilon]$  centered at  $y$ . Denoting by  $\hat{F}$  the distribution function of  $\hat{P}$ , the probability of that interval is

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

and its logarithm is equal to  $\ln \hat{f}(y)$  plus a constant that does not depend on  $\hat{f}$ . Minimizing (9) thus amounts to maximizing the probability of the observations.

In the case where the prediction takes the form of a RFS  $\tilde{Y}$ , we no longer have a single probability of the observation  $[y]_\epsilon$ , but two numbers: a degree of belief  $Bel_{\tilde{Y}}([y]_\epsilon)$  and a degree of plausibility  $Pl_{\tilde{Y}}([y]_\epsilon)$ . We could, thus, replace (9) by either

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

or

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

where the precision  $\epsilon$  now appears explicitly in the expression of the loss function. However, none of these two loss functions adequately measures the quality of the imprecise predictions as expressed by a RFS  $\tilde{Y}$ . To see this, let us assume that  $\tilde{Y} \sim \tilde{N}(\mu, \sigma^2, h)$  is a GRFN with precision  $h$ . From (7a), it is clear that  $Bel_{\tilde{Y}}([y]_\epsilon) \leq \Phi\left(\frac{y+\epsilon-\mu}{\sigma}\right) - \Phi\left(\frac{y-\epsilon-\mu}{\sigma}\right)$ , with an equality when  $h = +\infty$ . Consequently, whatever the values of  $\mu$  and  $\sigma^2$ , a Gaussian random variable with mean  $\mu$  and variance  $\sigma^2$  always achieves a loss smaller than that of  $\tilde{Y}$ , which means that a higher imprecision of the RFS is never rewarded by a decrease of the loss function. Now, when  $h = 0$ , we have  $Pl_{\tilde{Y}}([y]_\epsilon) = 1$  for any  $y, \mu$  and  $\sigma^2$ : the loss  $\underline{\mathcal{L}}_\epsilon(y, \tilde{Y})$  is, thus, maximized by the vacuous, uninformative RFS.

From the above analysis, we propose to consider as the loss function a weighted sum of  $\bar{\mathcal{L}}_\epsilon(y, \tilde{Y})$  and  $\underline{\mathcal{L}}_\epsilon(y, \tilde{Y})$ :

$$\mathcal{L}_{\lambda, \epsilon}(y, \tilde{Y}) = \lambda \bar{\mathcal{L}}_\epsilon(y, \tilde{Y}) + (1 - \lambda) \underline{\mathcal{L}}_\epsilon(y, \tilde{Y}), \quad (10)$$

where  $\lambda \in [0, 1]$  is a hyperparameter. Choosing a smaller value of  $\lambda$  amounts to favoring more cautious predictions. Figure 3 shows the variation of the loss  $\mathcal{L}_{\lambda, \epsilon}(0, \tilde{Y})$  with  $h$  for  $\lambda \in$

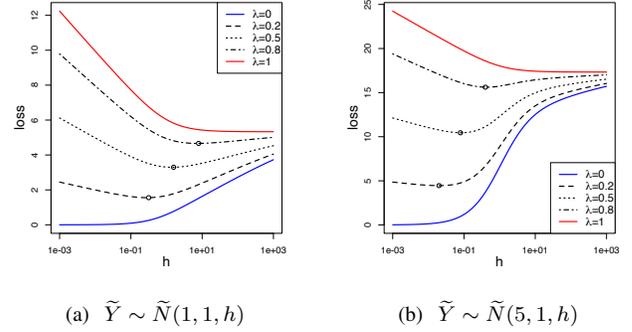


Fig. 3: Loss  $\mathcal{L}_{\lambda, \epsilon}(0, \tilde{Y})$  for a GRFN  $\tilde{Y} \sim \tilde{N}(\mu, 1, h)$  as a function of  $h$ , for  $\mu = 1$  (a) and  $\mu = 5$  (b) and  $\lambda \in \{0, 0.2, 0.5, 0.8, 1\}$ .

$\{0, 0.2, 0.5, 0.8, 1\}$ ,  $\epsilon = 0.01$ , and  $\tilde{Y} \sim \tilde{N}(\mu, 1, h)$  with  $\mu = 1$  (Figure 3a) and  $\mu = 5$  (Figure 3b). We can see that the loss is minimized for  $h = 0$  when  $\lambda = 0$ , for  $h \rightarrow +\infty$  when  $\lambda = 1$ , and for some intermediate value of  $h$  when  $0 < \lambda < 1$ . For any given value of  $\lambda$ , the optimum precision  $h$  is smaller when  $\mu = 5$ : a larger error is compensated by a smaller precision.

When  $\mu, \sigma^2$  and  $h$  are free to vary, loss function (10) is minimal for  $\mu \in [y]_\epsilon$ ,  $h = +\infty$  and  $\sigma^2 \rightarrow 0$ . As in support vector regression [28],  $\epsilon$  can be seen as a value under which the error is not penalized. We can also remark that, when  $\epsilon$  is small, we have, for a probabilistic GRFN  $\tilde{Y} \sim \tilde{N}(\mu, \sigma^2, +\infty)$ ,  $\mathcal{L}_{\lambda, \epsilon}(y, \tilde{Y}) \approx -\ln \phi\left(\frac{y-\mu}{\sigma}\right) - \ln \epsilon$ , where  $\phi$  denotes the standard normal probability density function; loss function (10) then approximates the negative log-likelihood, up to an additive constant.

In section IV, we will describe a neural network model for computing a GRFN  $\tilde{Y}(\mathbf{x}; \Psi)$  from an input vector  $\mathbf{x}$  and a parameter vector  $\Psi$ . Given a training set  $\mathcal{T} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ , the network will be trained by minimizing the empirical risk, or average loss

$$C_{\lambda, \epsilon}(\Psi) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\lambda, \epsilon}(y_i, \tilde{Y}(\mathbf{x}_i; \Psi)), \quad (11)$$

a finite-sample estimate of the expected loss  $\mathbb{E}_{\mathbf{X}, Y}[\mathcal{L}_{\lambda, \epsilon}(Y, \tilde{Y}(\mathbf{X}; \Psi))]$ .

### B. Calibration

In the case of probabilistic predictions of the response variable  $Y$ , we usually expect the predictive probability distribution to be close to the true conditional distribution of  $Y$  given  $\mathbf{X} = \mathbf{x}$ . More precisely, let  $\hat{F}_{Y|\mathbf{x}}$  be a predictive cumulative distribution function (cdf) for  $Y$  given  $\mathbf{X} = \mathbf{x}$ . For any  $\alpha \in (0, 1)$ , we can construct a prediction interval at level  $\alpha$  as

$$\mathcal{C}_\alpha(\mathbf{x}) = \left[ \hat{F}_{Y|\mathbf{x}}^{-1}\left(\frac{1-\alpha}{2}\right), \hat{F}_{Y|\mathbf{x}}^{-1}\left(\frac{1+\alpha}{2}\right) \right].$$

These intervals are said to be *calibrated* if

$$\forall \alpha \in (0, 1), \quad P_{\mathbf{X}, Y}(Y \in \mathcal{C}_\alpha(\mathbf{X})) \approx \alpha. \quad (12)$$

This property can be checked by estimating the coverage rates  $P_{\mathbf{X},Y}(Y \in \mathcal{C}_\alpha(\mathbf{X}))$  for different values of  $\alpha$  using a validation set and plotting the estimates vs. the nominal values  $\alpha$ . The resulting graph is known as a *calibration plot*.

In the case considered here, where the response is predicted for  $\mathbf{X} = \mathbf{x}$  by a GRFN  $\tilde{Y}(\mathbf{x})$ , we also need a notion of calibration that can easily be checked graphically. Different notions of calibration for belief functions are reviewed in [29]. Recently, Cella and Martin [30] proposed two definitions for “valid” predictions in the belief function framework. However, the weaker definition is too weak to be practically useful, while the stronger one appears to be verified only for consonant belief functions. Here, we will adopt a simple definition based on an immediate generalization of (12).

For any  $\alpha \in (0, 1]$ , we define an  $\alpha$ -level *belief prediction interval (BPI)* as an interval  $\mathcal{B}_\alpha(\mathbf{x}) = \mu(\mathbf{x}) \pm a(\mathbf{x})$  centered at  $\mu(\mathbf{x})$ , such that  $Bel_{\tilde{Y}(\mathbf{x})}(\mathcal{B}_\alpha(\mathbf{x})) = \alpha$ . A BPI at level  $\alpha$  is, thus, an interval that is believed to a degree  $\alpha$  to contain the true value of the response variable  $Y$ . The predictions will be said to be *calibrated* if, for all  $\alpha \in (0, 1]$ ,  $\alpha$ -level BPIs have a coverage probability at least equal to  $\alpha$ , i.e.,

$$\forall \alpha \in (0, 1], \quad P_{\mathbf{X},Y}(Y \in \mathcal{B}_\alpha(\mathbf{X})) \geq \alpha. \quad (13)$$

As in the probabilistic case, the calibration of evidential predictions can be checked graphically using a calibration plot (see Figure 6 in Section V-A). The predictions are calibrated if the curve lies above the first diagonal, and the predictions are all the more *precise* that the curve is close to the first diagonal.

#### IV. NEURAL NETWORK MODEL

In this section, we propose a neural network model, called ENNreg, which for an observed input vector  $\mathbf{X} = \mathbf{x}$  computes a GRFN  $\tilde{Y}(\mathbf{x})$ , with associated belief function  $Bel_{\tilde{Y}(\mathbf{x})}$  that “approximates”  $Y$  in some way. As the ENN model introduced in [9] for classification, ENNreg is based on prototypes. The distances to the prototypes are treated as independent pieces of evidence about the response and are combined by the generalized product-intersection rule  $\boxplus$  recalled in Section II-A.

The network architecture and propagation equations will first be described in Section IV-A. Learning will then be discussed in Section IV-B.

##### A. Neural network architecture

Let  $\mathbf{w}_1, \dots, \mathbf{w}_K$  denote  $K$  vectors in the  $p$ -dimensional feature space, called prototypes. The similarity between input vector  $\mathbf{x}$  and prototype  $\mathbf{w}_k$  is measured by

$$s_k(\mathbf{x}) = \exp(-\gamma_k^2 \|\mathbf{x} - \mathbf{w}_k\|^2), \quad (14)$$

where  $\gamma_k$  is a positive scale parameter. The evidence of prototype  $\mathbf{w}_k$  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  $\sigma_k^2$  and  $h_k$  are variance and precision parameters for prototype  $k$ ; the mean  $\mu_k(\mathbf{x})$  is defined as

$$\mu_k(\mathbf{x}) = \beta_k^T \mathbf{x} + \beta_{k0}, \quad (15)$$

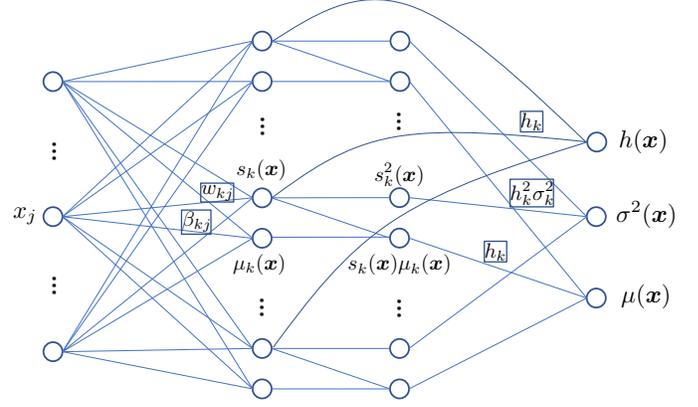


Fig. 4: Neural network architecture.

where  $\beta_k$  is a  $p$ -dimensional vector of coefficients, and  $\beta_{k0}$  is a scalar parameter. The quantity  $\mu_k(\mathbf{x})$  can be seen as an estimate of the conditional expectation of the response  $Y$  given that  $\mathbf{x}$  is close to  $\mathbf{w}_k$ ,  $\sigma_k^2$  is an estimate of the conditional variance, while the precision term  $s_k(\mathbf{x})h_k$  represents epistemic uncertainty. Obviously, when the distance  $\|\mathbf{x} - \mathbf{w}_k\|$  tends to infinity, the precision  $s_k(\mathbf{x})h_k$  tends to 0 and  $\tilde{Y}_k(\mathbf{x})$  becomes vacuous. The vector  $\psi_k$  of parameters associated to prototype  $k$  is, thus,  $\psi_k = (\mathbf{w}_k, \gamma_k, \beta_k, \beta_{k0}, \sigma_k^2, h_k)$ .

The output  $\tilde{Y}(\mathbf{x})$  for input  $\mathbf{x}$  is computed by combining the GRFNs  $\tilde{Y}_k(\mathbf{x})$ ,  $k = 1, \dots, K$  induced by the  $K$  prototypes using the  $\boxplus$  operator. From (8), it is a GRFN  $\tilde{Y}(\mathbf{x}) \sim \tilde{N}(\mu(\mathbf{x}), \sigma^2(\mathbf{x}), h(\mathbf{x}))$ , with

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

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

and  $h(\mathbf{x}) = \sum_{k=1}^K s_k(\mathbf{x})h_k$ .

These operations can be seen as being performed by a neural network with an input layer, two hidden layers of  $2K$  units and an output layer of three units, as shown in Figure 4. The first hidden layer is composed of  $K$  pairs of units totally connected to the input units: a radial basis function (RBF) unit with weight vector  $\mathbf{w}_k$  that computes  $s_k(\mathbf{x})$  and a linear unit with weight vector  $\beta_k$  and bias  $\beta_{k0}$  that computes  $\mu_k(\mathbf{x})$ . The second layer is also composed of  $K$  pairs of units: in each pair  $k$ , a “squaring” unit connected to RBF unit  $k$  of the previous layer computes  $s_k^2(\mathbf{x})$ , and a “product” unit connected to RBF unit  $k$  and linear unit  $k$  of the previous layer computes the product  $s_k(\mathbf{x})\mu_k(\mathbf{x})$ . Finally, the output layer is composed of (1) a linear unit with shortcut connections from the RBF units in the first hidden layer and weights ( $h_k$ ) that computes  $h(\mathbf{x})$ ; (2) a unit connected to the squaring units of the previous layer with weights ( $h_k^2\sigma_k^2$ ) that outputs  $\sigma^2(\mathbf{x})$ , and (3) a unit connected to the product units of the previous layer with weights ( $h_k$ ) that computes  $\mu(\mathbf{x})$ .

Although the structure of this network is more complex than that of RBF networks [31] [32], the complexity of one input

propagation is the same, i.e.,  $O(pK)$ . We can also remark that this neural network model generalizes both RBF networks and the linear model:

- If  $\beta_k = 0$  for all  $k$ , then  $\mu_k(\mathbf{x}) = \beta_{k0}$ , and  $\mu(\mathbf{x})$  is identical to the output of an RBF neural network [33] with  $K$  hidden units, hidden-to-output weights  $h_k\beta_{k0}$  and normalized outputs;
- If  $\gamma_k = 0$  for all  $k$ , the output  $\mu(\mathbf{x})$  becomes linear in  $\mathbf{x}$  and the variance component  $\sigma^2(\mathbf{x})$  is constant. We then have a linear model with constant variance.

### B. Learning

Let  $\Psi = (\psi_1, \dots, \psi_K)$  denote the vector of all parameters in the model. Given a training set  $\mathcal{T} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ , the accuracy of the predictions can be measured by the empirical risk (11) introduced in Section III-A. To train the ENNreg model described in Section IV-A, we will use the regularized average loss

$$C_{\lambda, \epsilon, \xi, \rho}^{(R)}(\Psi) = C_{\lambda, \epsilon}(\Psi) + \frac{\xi}{K} \sum_{k=1}^K h_k + \frac{\rho}{K} \sum_{k=1}^K \gamma_k^2, \quad (16)$$

where  $\xi$  and  $\rho$  are regularization coefficients. The first regularization term has the effect of reducing the influence of some prototypes (setting  $h_k = 0$  actually amounts to discarding prototype  $k$ ), while the second one shrinks the solution towards a linear model. Loss function (16) can be minimized using any batch or online gradient-based optimization procedure. The constraints  $h_k \geq 0$  are imposed by introducing intermediate variables  $\eta_k$  such that  $h_k = \eta_k^2$ . The gradient can be computed by error back-propagation, with complexity  $O(Kp)$ .

As in RBF networks, the prototypes in the ENNreg model can easily be initialized using a clustering procedure such as the  $K$ -means algorithm. As before, let  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  be the training set. Let  $\mathcal{I}_k$  denote the set of indices of input vectors in cluster  $k$ , and  $n_k = |\mathcal{I}_k|$ . The prototypes are initialized as the cluster centers,

$$\mathbf{w}_k = \frac{1}{n_k} \sum_{i \in \mathcal{I}_k} \mathbf{x}_i,$$

and  $\gamma_k$  is initialized proportionally to the inverse square root of the mean squared distances within cluster  $k$ :

$$\gamma_k = \frac{1}{\sqrt{2}} \left( \sum_{i \in \mathcal{I}_k} \frac{1}{n_k} \|\mathbf{x}_i - \mathbf{w}_k\|^2 \right)^{-1/2}.$$

Parameters  $\beta_{k0}$  and  $\sigma_k^2$  are set, respectively, to the mean and variance of the response variable within cluster  $k$ :

$$\beta_{k0} = \frac{1}{n_k} \sum_{i \in \mathcal{I}_k} y_i, \quad \sigma_k^2 = \frac{1}{n_k} \sum_{i \in \mathcal{I}_k} (y_i - \beta_{k0})^2.$$

Finally, we set  $\beta_k = 0$  and  $h_k = 1$ .

Overall, the model has five hyperparameters: the number  $K$  of prototypes, coefficients  $\epsilon$  and  $\lambda$  in the definition of the loss (10), and regularization coefficients  $\xi$  and  $\rho$  in (16). In practice,  $K$  can be set to a large value, and the effective number of prototypes can be controlled by  $\xi$ . We found the network performance to be quite robust to the choice of  $\epsilon$  and  $\lambda$ . In

the experiments reported in Section V, these hyperparameters were fixed to  $\epsilon = 0.01\hat{\sigma}_Y$ , where  $\hat{\sigma}_Y$  is the sample standard deviation of  $Y$ , and  $\lambda = 0.9$ . Only  $\xi$  and  $\rho$  thus have to be tuned using either cross-validation or the hold-out method.

## V. EXPERIMENTAL RESULTS

We first give an illustrative example in Section V-A. Results from a comparative experiment are then reported in Section V-B.

### A. Illustrative example

As an illustrative example, we consider independent and identically distributed (iid) data simulated from the following distribution: the one-dimensional input  $X$  has a uniform distribution in the interval  $[-2, 2]$ , and

$$Y = X + (\sin 3X)^3 + \frac{X+2}{4\sqrt{2}}U,$$

where  $U \sim N(0, 1)$  is a standard normal variable. The conditional standard deviation of  $Y$  given  $X = x$  thus increases linearly from 0 for  $x = -2$  to  $1/\sqrt{2}$  for  $x = 2$ .

We generated a learning set of size  $n = 300$  and a validation set of the same size. Figure 5a shows the learning data with the expected values  $\mu(x)$ , as well as the intervals  $\mu(x) \pm 2\sigma(x)$  as functions of  $x$  for a network initialized with  $K = 30$  prototypes, before training. We can see that the initial solution is already a good approximation to the true regression function.

The values of hyperparameters  $\xi$  and  $\rho$  minimizing the mean squared error on the validation set were  $\xi = 10^{-4}$  and  $\rho = 0.01$ . Figure 5b shows the outputs  $\mu(x)$  of a network trained with these values, together with BPIs at levels  $\alpha \in \{0.5, 0.9, 0.99\}$ .

A calibration plot is shown in Figure 6, representing the proportion of  $\alpha$ -level BPIs containing the observed value of the response in the validation set, for  $\alpha \in \{0.1, 0.2, \dots, 0.9\}$ . In Figure 6, we also show the calibration graph for the probabilistic prediction intervals defined as  $\mu(x) \pm u_{(1+\alpha)/2}\sigma(x)$ . In this particular case, both prediction intervals are calibrated, the BPIs being more conservative than the probabilistic ones. We can remark that we can make the BPIs less conservative by multiplying the output precisions  $h(x)$  by some constant  $c > 1$ . The probabilistic intervals are recovered when  $c \rightarrow +\infty$ . Conversely, we can make the BPIs more conservative by multiplying the output precision by a constant  $c < 1$ . This technique will be illustrated in Section V-B.

From Figure 6, we can see that the calibration curve of BPIs is close to that of probabilistic prediction intervals, which indicates that, for this dataset, the output precision  $h(x)$  is quite high. This is confirmed by Figure 7a showing  $h(x)$  vs.  $x$ . Precision first increases, and then decreases with  $x$  as the errors  $|y_i - \mu(x_i)|$  increase, which can be explained by Figure 3: larger errors result in more imprecise predictions. Figure 7b shows the output standard deviation  $\sigma(x)$  as a function of  $x$ : it slightly overestimates the true standard deviation, but the increasing trend is well captured.

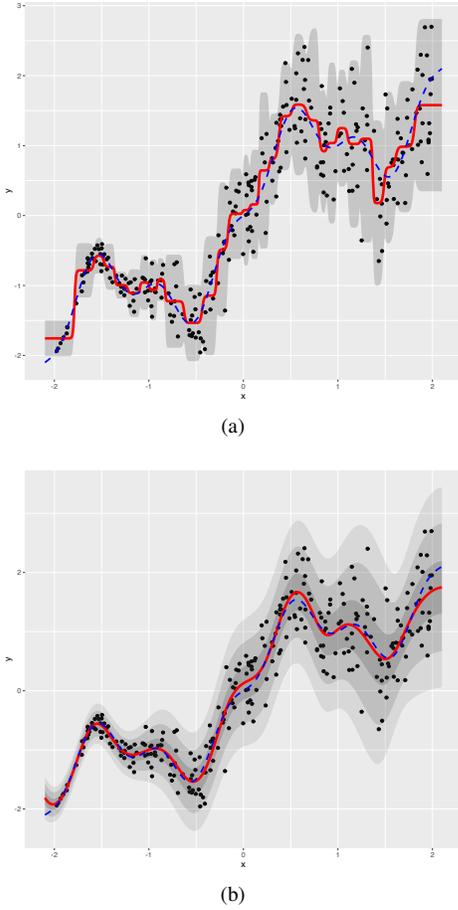


Fig. 5: (a) Learning data, true regression function (blue broken line) and expected values  $\mu(x)$  with a neural network of  $K = 30$  prototypes after initialization (red solid line) with a two-standard-deviation intervals (grey area); (b) Predictions after training: expected values  $\mu(x)$  (red solid line) and belief prediction intervals at levels  $\alpha \in \{0.5, 0.9, 0.99\}$  (grey areas).

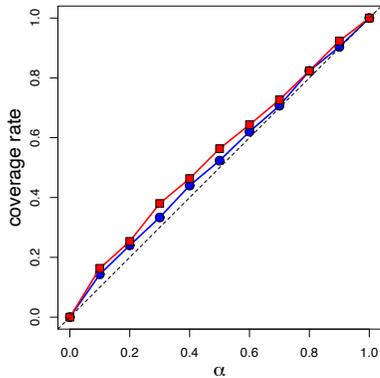


Fig. 6: Calibration plots for probabilistic prediction intervals (blue circles) and belief prediction intervals (red squares), showing the validation coverage rates vs. the belief or confidence levels  $\alpha \in \{0.1, \dots, 0.9\}$ .

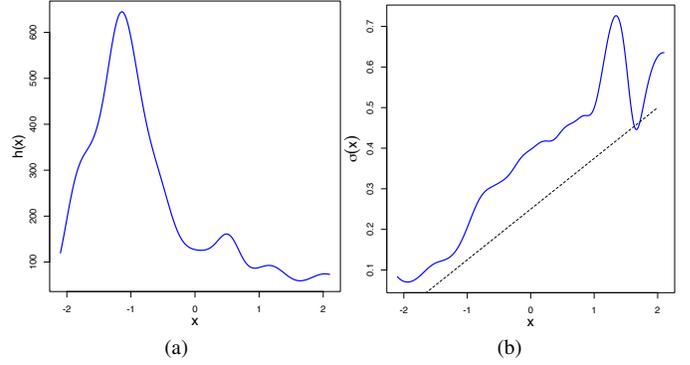


Fig. 7: Output precision  $h(x)$  (a) and standard deviation  $\sigma(x)$  (b) for the simulated data. The broken line indicates the true standard deviation.

TABLE I: Characteristics of the datasets.

	$n$	$p$	response
Boston	506	13	medv
Energy	768	8	Y2
Concrete	1030	8	strength
Yacht	308	6	Y
Wine	1599	11	quality
kin8nm	8192	8	V9
Crime	1994	100	ViolentCrimesPerPop
Residential	372	103	V10
Airfoil	1503	5	Y
Bike	731	9	cnt

### B. Comparative experiments

We compared the performance of ENNreg in terms of root mean squared (RMS) error to those of eight alternative regression methods on ten datasets. Eight of the datasets (Energy efficiency, Concrete Compressive Strength, Yacht Hydrodynamics, Wine quality – red, Communities and Crime, Residential building, Airfoil Self-Noise, Bike sharing) are available from the UCI Machine Learning Repository<sup>3</sup>. The Boston dataset is included in the R package MASS [34], and the kin8nm dataset was downloaded from the OpenML web site<sup>4</sup>. The characteristics of these datasets are summarized in Table I. We also give in Table I the name of the response variable, as some datasets have several possible response variables. For the “Communities and Crime” dataset, variables with two missing values or more were discarded. For the “Bike sharing dataset”, the nine predictors were season, year, month, holiday, weekday, weathersit, atemp, hum and windspeed.

The eight methods are:

- Four nonlinear regression algorithms using RBF kernel functions: RBF networks, Relevance Vector Machines (RVM), Support Vector Machines (SVM), and Gaussian Processes (GP);
- Two state-of-the-art methods for nonlinear regression: Random Forests (RF) and Multilayer Perceptrons (MLP) with one hidden layer of sigmoidal units;

<sup>3</sup><https://archive.ics.uci.edu/ml/>.

<sup>4</sup><https://www.openml.org>.

TABLE II: Average RMS and standard errors for ENNreg and eight alternative algorithms on ten datasets.

	ENNreg	RBF	RVM	SVM	GP	RF	MLP	ridge	lasso
Boston	<b>2.87 ± 0.14</b>	3.31 ± 0.19	3.42 ± 0.17	3.17 ± 0.15	3.70 ± 0.22	<b>3.11 ± 0.14</b>	<b>3.14 ± 0.14</b>	5.05 ± 0.23	5.02 ± 0.21
Energy	<b>1.06 ± 0.05</b>	2.06 ± 0.08	1.79 ± 0.05	1.39 ± 0.06	2.58 ± 0.07	1.75 ± 0.06	<b>0.95 ± 0.16</b>	3.56 ± 0.10	3.26 ± 0.09
Concr.	5.10 ± 0.12	6.30 ± 0.19	6.38 ± 0.16	5.62 ± 0.13	6.93 ± 0.13	<b>4.64 ± 0.12</b>	<b>4.82 ± 0.16</b>	10.71 ± 0.17	10.63 ± 0.18
Yacht	<b>0.44 ± 0.04</b>	2.00 ± 0.20	1.88 ± 0.20	1.93 ± 0.11	6.12 ± 0.31	0.96 ± 0.08	<b>0.50 ± 0.05</b>	8.48 ± 0.24	8.35 ± 0.23
Wine	0.63 ± 0.01	0.63 ± 0.01	0.80 ± 0.02	0.61 ± 0.01	0.61 ± 0.01	<b>0.56 ± 0.01</b>	0.77 ± 0.01	0.65 ± 0.01	0.65 ± 0.01
kin8nm	0.08 ± 0.00	0.11 ± 0.00	–	0.09 ± 0.00	0.08 ± 0.00	0.14 ± 0.00	<b>0.07 ± 0.00</b>	0.20 ± 0.00	0.20 ± 0.00
Crime	<b>0.14 ± 0.00</b>								
Resid.	<b>0.11 ± 0.01</b>	0.16 ± 0.01	0.17 ± 0.01	0.15 ± 0.01	0.22 ± 0.01	0.16 ± 0.01	0.14 ± 0.01	0.18 ± 0.01	0.17 ± 0.01
Airfoil	<b>1.46 ± 0.03</b>	1.70 ± 0.04	2.58 ± 0.04	2.37 ± 0.04	2.49 ± 0.04	<b>1.44 ± 0.04</b>	1.53 ± 0.04	3.84 ± 0.04	3.83 ± 0.04
Bike	<b>6.59 ± 0.19</b>	<b>6.49 ± 0.15</b>	<b>6.64 ± 0.14</b>	7.11 ± 0.16	7.55 ± 0.14	6.86 ± 0.17	9.68 ± 0.20	7.82 ± 0.16	7.76 ± 0.17

TABLE III: Average CPU time and standard errors for ENNreg and RBF networks on five datasets.

Dataset	ENNreg	RBF
Boston	2.63 ± 0.09	2.88 ± 0.08
Concrete	2.45 ± 0.08	8.16 ± 0.19
Wine	1.92 ± 0.03	4.97 ± 0.24
Resid.	11.38 ± 0.47	23.43 ± 0.55
Airfoil	20.50 ± 1.30	28.90 ± 0.39

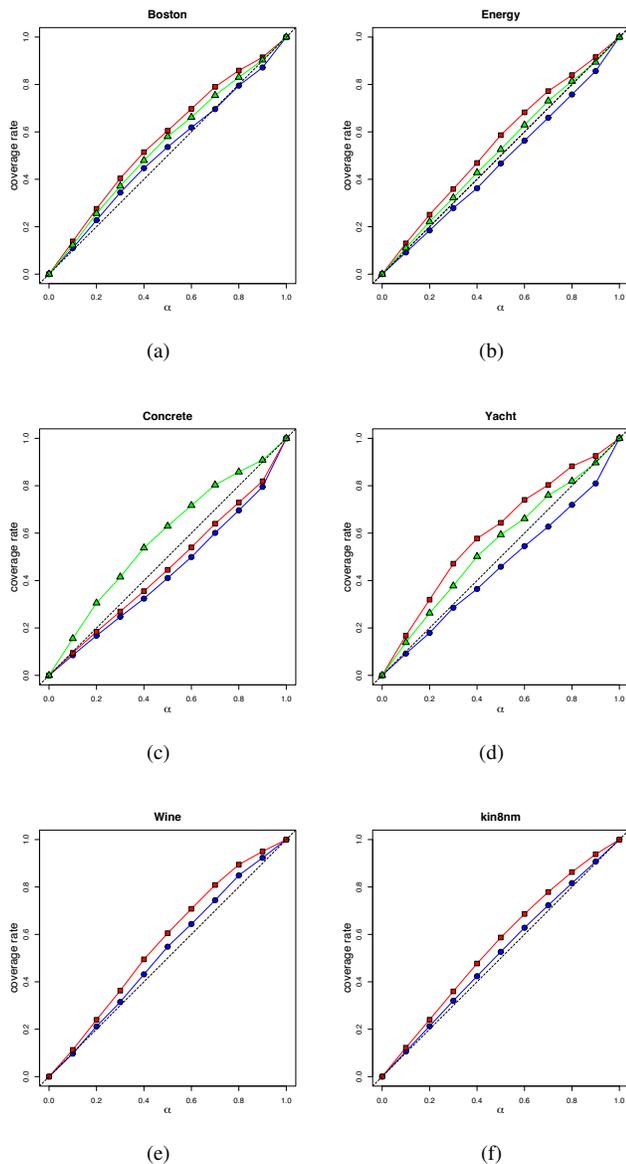


Fig. 8: Calibration plots for ENNreg and six datasets: probabilistic predictions (blue circles), raw evidential predictions (red squares) and adjusted evidential predictions (green triangles).

- Two regularized linear regression methods: Ridge and Lasso.

For all methods, except RBF networks, we used the implementation in the R package `caret` [35]. For training ENNreg, we used batch learning using the algorithm described in [36] for all datasets, except `kin8nm` for which we used mini-batch stochastic gradient descent. Each dataset was split randomly into training and test sets containing, respectively, 90% and 10% of the observations. The random splits were repeated 20 times. All input variables were scaled to have zero mean and unit standard deviation. For each method, hyperparameters were tuned by 5-fold cross-validation. For ENNreg, we set  $\lambda = 0.9$  and  $\epsilon = 0.01\sigma_y$  (where  $\sigma_y$  is the standard deviation of the response variable) for all the simulations. The number of prototypes was fixed to  $K = 30$  for all datasets except `kin8nm` and `Airfoil`, for which it was set to  $K = 100$ . Hyperparameters  $\xi$  and  $\rho$  was tuned by cross-validation.

The results<sup>5</sup> are reported in Table II. We can see that ENNreg stands among the best methods for seven out of the ten datasets. It is only outperformed by RF and MLP on the Concrete dataset, RF on the Wine dataset, and MLP on the `kin8nm` dataset. In most cases, it performs strictly better than other RBF-based methods (except for the Crime dataset, on which all methods have equivalent performances, and for the Bike dataset, on which it performs as well as RBF and RVM). Overall, it appears that ENNreg is a very competitive regression method in terms of prediction accuracy.

Comparing computing time across all methods is difficult because it highly depends on implementation. We have seen in Section IV that both forward and back propagation in ENNreg can be performed in linear time with respect to both the number of inputs and the number of prototypes; consequently,

<sup>5</sup>Results with the RVM algorithm are not reported for the `kin8nm` dataset, because of the excessive computing time on our machine (a 2019 16" MacBook Pro with a 2.4 GHz 8-core Intel i9 processor).

TABLE IV: Average RMS and standard errors for ENNreg, probabilistic backpropagation (PBP), Monte Carlo dropout and deep ensembles for six datasets.

	ENNreg	PRP	MC-dropout	Deep ens.
Boston	<b>2.87 ± 0.14</b>	<b>3.01 ± 0.18</b>	<b>2.97 ± 0.19</b>	<b>3.28 ± 1.00</b>
Energy	<b>1.06 ± 0.05</b>	1.80 ± 0.05	1.66 ± 0.04	2.09 ± 0.29
Concr.	<b>5.10 ± 0.12</b>	5.67 ± 0.09	<b>5.23 ± 0.12</b>	6.03 ± 0.58
Yacht	<b>0.44 ± 0.04</b>	1.02 ± 0.05	1.11 ± 0.09	1.58 ± 0.48
Wine	<b>0.63 ± 0.01</b>	<b>0.64 ± 0.01</b>	<b>0.62 ± 0.01</b>	<b>0.64 ± 0.04</b>
kin8nm	<b>0.08 ± 0.00</b>	0.10 ± 0.00	0.10 ± 0.00	0.09 ± 0.00

TABLE V: Average NLL and standard errors for ENNreg, probabilistic backpropagation (PBP), Monte Carlo dropout and deep ensembles for six datasets.

	ENNreg	PBP	MC-dropout	Deep ens.
Boston	<b>2.53 ± 0.07</b>	<b>2.57 ± 0.09</b>	<b>2.46 ± 0.06</b>	<b>2.41 ± 0.25</b>
Energy	<b>1.14 ± 0.07</b>	2.04 ± 0.02	1.99 ± 0.02	<b>1.38 ± 0.22</b>
Concr.	3.38 ± 0.13	3.16 ± 0.02	<b>3.04 ± 0.02</b>	<b>3.06 ± 0.18</b>
Yacht	<b>0.13 ± 0.12</b>	1.63 ± 0.02	1.55 ± 0.03	1.18 ± 0.21
Wine	<b>0.94 ± 0.01</b>	0.97 ± 0.01	<b>0.93 ± 0.01</b>	<b>0.94 ± 0.12</b>
kin8nm	<b>-1.19 ± 0.00</b>	-0.90 ± 0.01	-0.95 ± 0.01	<b>-1.20 ± 0.02</b>

the algorithmic efficiency of ENNreg is similar to those of other neural network methods. We have done a comparison with the RBF network model because it was implemented in a similar way as ENNreg and with the same optimization algorithm. Average computing times for both methods on five datasets are shown in Table III. We can see that ENNreg is actually faster than RBF networks with the same number of prototypes, which can be explained by faster convergence due to greater flexibility.

We also compared our results to published results obtained with three recent neural network methods: probabilistic backpropagation (PBP) [37], Monte Carlo dropout [38], and deep ensembles [39], on six of the ten datasets for which published results with these methods are available. The RMS values are shown in Table IV. We can see that ENNreg performs strictly better than, or as well as these three methods on the six considered datasets in terms of RMS.

The PBP, Monte Carlo dropout and deep ensemble methods were actually introduced for uncertainty quantification in deep networks. Probabilistic calibration is usually measured by the negative loglikelihood (NLL), assuming Gaussian errors. For ENNreg, NLL can be computed using the mean  $\mu(\mathbf{x})$  and variance  $\sigma^2(\mathbf{x})$  output, ignoring the precision  $h(\mathbf{x})$ . ENNreg is then seen as a probabilistic method. The NLL values for ENNreg and the three neural network methods on the six datasets (as reported in [37], [38] and [39]) are shown in Table V. We can see that ENNreg stands among the best methods for five out of the six datasets in terms of probabilistic calibration as measured by NLL.

As mentioned in Section V-A, the ENNreg model provides not only conditional mean and variance estimates as probabilistic models, but a precision output  $h(\mathbf{x})$  that can be used to compute calibrated BPIs as defined by (13). Figure 8 shows calibration plots for the same six datasets as above. Each plot shows the calibration curves of (1) probabilistic prediction intervals (using only  $\mu(\mathbf{x})$  and  $\sigma^2(\mathbf{x})$ ), (2) raw BPIs and (3)

adjusted BPIs obtained by multiplying the output  $h(\mathbf{x})$  by a constant  $c > 0$  in such a way that the calibration curve lies above the diagonal, but as close as possible to it. We can distinguish several cases. For the Boston, Energy and Yacht datasets (Figures 8a, 8b and 8d), the BPIs are calibrated but too conservative. Applying a correction factor  $c = 2$  increased the precision of the predictions. In the case of the Concrete dataset (Figure 8c), the BPIs before adjustment are not calibrated. In this case, we need to apply a correction factor  $c < 1$  to make the predictions more imprecise. Here, we used  $c = 0.1$ . Finally, for the Wine and kin8nm datasets (Figures 8e and 8f), the calibration curves corresponding to probabilistic predictions are already above the diagonal, which means that the probabilistic prediction intervals are already calibrated BPIs. In such a case, we can consider that the probabilistic predictions provide an adequate description of prediction uncertainty and neglect the precision output  $h(\mathbf{x})$  or, equivalently, multiply it by  $c = \infty$ .

## VI. CONCLUSIONS

We have introduced an evidential neural network model for regression, called ENNreg. In this model, distances to  $K$  prototypes are considered as pieces of evidence about the response variable and are described by GRFNs. The total evidence is then pooled by the generalized product-intersection rule, an extension of Dempster’s rule in the epistemic random fuzzy set setting. The network architecture is composed of a first hidden layer of  $K$  RBF units and  $K$  linear units, a second hidden layer of  $2K$  units, and an output layer of three units with cross-cut connections from the two hidden units. The network output is a GRFN described by three numbers: a point prediction  $\mu(\mathbf{x})$ , a conditional variance estimate  $\sigma^2(\mathbf{x})$ , and a precision parameter  $h(\mathbf{x})$  whose value depends on the distances between the input vector and the prototypes, and which can be seen as describing the reliability of the probabilistic predictions. This additional degree of freedom makes it possible to quantify not only random uncertainty, but also epistemic uncertainty.

The network output can also be expressed as a predictive belief function on the real line induced by the output GRFN. We have defined a loss function for such outputs, extending the negative log-likelihood to evidential regression. We have also discussed the calibration of predictive belief functions for regression tasks and introduced a definition based on the coverage rates of belief prediction intervals (prediction intervals with specified belief degree), which are required to be conservative. Calibrated belief prediction intervals at degree  $\alpha$  thus contain, on average, at least  $100\alpha\%$  of future observations. This calibration property can be assessed visually by drawing calibration plots.

Comparative experiments have shown that ENNreg performs better in terms of RMS error than other RBF-based regression methods such as RBF networks, SVM, RVM and Gaussian processes, and to be competitive with state-of-the-art methods for regression such as random forests and multilayer perceptrons. Furthermore, a comparison with results reported in the literature for probabilistic back-propagation, Monte

Carlo dropout and deep ensemble shows that ENNreg is also competitive with these deep learning approaches in terms of both prediction accuracy and probabilistic calibration, while offering the advantage of a richer quantification of uncertainty thanks to the random fuzzy set formalism.

This research can be extended in several directions. For structured inputs such as time series or images, additional feature-extraction layers could be inserted before the prototype layer, as done in [10] for image classification. For regression with multiple outputs, the ENNreg model could be extended to compute outputs in the form of Gaussian random fuzzy vectors as introduced in [19]. Finally, other families of random fuzzy numbers could be used to accommodate outputs with, e.g., skewed distribution or bounded support.

## REFERENCES

- [1] A. P. Dempster, "Upper and lower probabilities induced by a multivalued mapping," *Annals of Mathematical Statistics*, vol. 38, pp. 325–339, 1967.
- [2] G. Shafer, *A mathematical theory of evidence*. Princeton, N.J.: Princeton University Press, 1976.
- [3] T. Denœux, D. Dubois, and H. Prade, "Representations of uncertainty in artificial intelligence: Beyond probability and possibility," in *A Guided Tour of Artificial Intelligence Research*, P. Marquis, O. Papini, and H. Prade, Eds. Springer Verlag, 2020, vol. 1, ch. 4, pp. 119–150.
- [4] T. Denœux, S. Li, and S. Sriboonchitta, "Evaluating and comparing soft partitions: an approach based on Dempster-Shafer theory," *IEEE Transactions on Fuzzy Systems*, vol. 26, no. 3, pp. 1231–1244, 2018.
- [5] Z. Su and T. Denœux, "BPEC: Belief-peaks evidential clustering," *IEEE Transactions on Fuzzy Systems*, vol. 27, no. 1, pp. 111–123, 2019.
- [6] T. Denœux, "A  $k$ -nearest neighbor classification rule based on Dempster-Shafer theory," *IEEE Trans. on Systems, Man and Cybernetics*, vol. 25, no. 05, pp. 804–813, 1995.
- [7] Z.-G. Liu, Y. Liu, J. Dezert, and F. Cuzzolin, "Evidence combination based on credal belief redistribution for pattern classification," *IEEE Transactions on Fuzzy Systems*, vol. 28, no. 4, pp. 618–631, 2020.
- [8] Z. Liu, X. Zhang, J. Niu, and J. Dezert, "Combination of classifiers with different frames of discernment based on belief functions," *IEEE Transactions on Fuzzy Systems*, vol. 29, no. 7, pp. 1764–1774, 2021.
- [9] T. Denœux, "A neural network classifier based on Dempster-Shafer theory," *IEEE Trans. on Systems, Man and Cybernetics A*, vol. 30, no. 2, pp. 131–150, 2000.
- [10] Z. Tong, P. Xu, and T. Denœux, "An evidential classifier based on Dempster-Shafer theory and deep learning," *Neurocomputing*, vol. 450, pp. 275–293, 2021.
- [11] 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*, vol. 149, pp. 39–60, 2022.
- [12] A. P. Dempster, "Upper and lower probabilities generated by a random closed interval," *Annals of Mathematical Statistics*, vol. 39, no. 3, pp. 957–966, 1968.
- [13] P. Smets, "Belief functions on real numbers," *International Journal of Approximate Reasoning*, vol. 40, no. 3, pp. 181–223, 2005.
- [14] O. Kanjanatarakul, S. Sriboonchitta, and T. Denœux, "Prediction of future observations using belief functions: A likelihood-based approach," *International Journal of Approximate Reasoning*, vol. 72, pp. 71–94, 2016.
- [15] T. Denœux, "Function approximation in the framework of evidence theory: A connectionist approach," in *Proceedings of the 1997 International Conference on Neural Networks (ICNN'97)*, vol. 1. Houston: IEEE, June 1997, pp. 199–203.
- [16] S. Petit-Renaud and T. Denœux, "Handling different forms of uncertainty in regression analysis: a fuzzy belief structure approach," in *Symbolic and quantitative approaches to reasoning and uncertainty (ECSQARU'99)*, A. Hunter and S. Pearsons, Eds. London: Springer Verlag, June 1999, pp. 340–351.
- [17] —, "Nonparametric regression analysis of uncertain and imprecise data using belief functions," *International Journal of Approximate Reasoning*, vol. 35, no. 1, pp. 1–28, 2004.
- [18] T. Denœux, "An evidential neural network model for regression based on random fuzzy numbers," in *Belief Functions: Theory and Applications*, S. Le Hégarat-Masclé, I. Bloch, and E. Aldea, Eds. Cham: Springer International Publishing, 2022, pp. 57–66.
- [19] T. Denœux, "Reasoning with fuzzy and uncertain evidence using epistemic random fuzzy sets: General framework and practical models," *Fuzzy Sets and Systems*, 2022, <https://doi.org/10.1016/j.fss.2022.06.004>.
- [20] T. Denœux, "Belief functions induced by random fuzzy sets: A general framework for representing uncertain and fuzzy evidence," *Fuzzy Sets and Systems*, vol. 424, pp. 63–91, 2021.
- [21] H. T. Nguyen, "On random sets and belief functions," *Journal of Mathematical Analysis and Applications*, vol. 65, pp. 531–542, 1978.
- [22] M. L. Puri and D. A. Ralescu, "Fuzzy random variables," *Journal of Mathematical Analysis and Applications*, vol. 114, no. 2, pp. 409–422, 1986.
- [23] S. de la Rosa de Sáa, M. A. Lubiano, B. Sinova, M. A. Gil, and P. Filzmoser, "Location-free robust scale estimates for fuzzy data," *IEEE Transactions on Fuzzy Systems*, vol. 29, no. 6, pp. 1682–1694, 2021.
- [24] J. C. Figueroa-Garca, C. A. Varn-Gaviria, and J. L. Barbosa-Fontecha, "Fuzzy random variable generation using  $\alpha$ -cuts," *IEEE Transactions on Fuzzy Systems*, vol. 29, no. 3, pp. 539–548, 2021.
- [25] I. Couso and L. Sánchez, "Upper and lower probabilities induced by a fuzzy random variable," *Fuzzy Sets and Systems*, vol. 165, no. 1, pp. 1–23, 2011.
- [26] L. A. Zadeh, "Fuzzy sets as a basis for a theory of possibility," *Fuzzy Sets and Systems*, vol. 1, pp. 3–28, 1978.
- [27] D. Dubois, H. Prade, and R. Yager, "Merging fuzzy information," in *Fuzzy sets in approximate reasoning and information systems*, J. C. Bezdek, D. Dubois, and H. Prade, Eds. Boston: Kluwer Academic Publishers, 1999, pp. 335–401.
- [28] H. Drucker, C. J. C. Burges, L. Kaufman, A. Smola, and V. Vapnik, "Support vector regression machines," in *Advances in Neural Information Processing Systems*, M. Mozer, M. Jordan, and T. Petsche, Eds., vol. 9. MIT Press, 1996, pp. 155–161.
- [29] T. Denœux and S. Li, "Frequency-calibrated belief functions: Review and new insights," *International Journal of Approximate Reasoning*, vol. 92, pp. 232–254, 2018.
- [30] L. Cella and R. Martin, "Valid inferential models for prediction in supervised learning problems," *International Journal of Approximate Reasoning*, vol. 150, pp. 1–18, 2022.
- [31] J. Moody and C. J. Darken, "Fast learning in networks of locally-tuned processing units," *Neural Computation*, vol. 1, no. 2, pp. 281–294, 1989.
- [32] M. Musavi, W. Ahmed, K. Chan, K. Faris, and D. Hummels, "On the training of radial basis function classifiers," *Neural Networks*, vol. 5, no. 4, pp. 595–603, 1992.
- [33] D. S. Broomhead and D. Lowe, "Multivariable functional interpolation and adaptive networks," *Complex Systems*, vol. 2, no. 3, pp. 321–355, 1988.
- [34] W. N. Venables and B. D. Ripley, *Modern Applied Statistics with S*, 4th ed. New York: Springer, 2002.
- [35] M. Kuhn, *caret: Classification and Regression Training*, 2021, r package version 6.0-90, <https://CRAN.R-project.org/package=caret>. [Online]. Available: <https://CRAN.R-project.org/package=caret>
- [36] F. M. Silva and L. B. Almeida, "Speeding up backpropagation," in *Advanced neural computers*, R. Eckmiller, Ed. New-York: Elsevier-North-Holland, 1990, pp. 151–158.
- [37] J. M. Hernandez-Lobato and R. Adams, "Probabilistic backpropagation for scalable learning of Bayesian neural networks," in *Proceedings of the 32nd International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, F. Bach and D. Blei, Eds., vol. 37. Lille, France: PMLR, 07–09 Jul 2015, pp. 1861–1869.
- [38] Y. Gal and Z. Ghahramani, "Dropout as a Bayesian approximation: Representing model uncertainty in deep learning," in *Proceedings of The 33rd International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, M. F. Balcan and K. Q. Weinberger, Eds., vol. 48. New York, New York, USA: PMLR, 20–22 Jun 2016, pp. 1050–1059.
- [39] B. Lakshminarayanan, A. Pritzel, and C. Blundell, "Simple and scalable predictive uncertainty estimation using deep ensembles," in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30. Curran Associates, Inc., 2017.



**Thierry Denœux** Thierry Denœux is a Full Professor (Exceptional Class) with the Department of Information Processing Engineering at the University of Compiègne, France, and a senior member of the French Academic Institute (*Institut Universitaire de France*). His research interests concern reasoning and decision-making under uncertainty and, more generally, the management of uncertainty in intelligent systems. His main contributions are in the theory of belief functions with applications to statistical inference, pattern recognition, machine learning and information fusion. He has published more than 300 papers in this area. He is the Editor-in-Chief of the *International Journal of Approximate Reasoning*, and an Associate Editor of several journals including *Fuzzy Sets and Systems* and *International Journal on Uncertainty, Fuzziness and Knowledge-Based Systems*.