



Learning Sets of Probabilities Through Ensemble Methods

Vu-Linh Nguyen, Haifei Zhang, Sébastien Destercke

► To cite this version:

Vu-Linh Nguyen, Haifei Zhang, Sébastien Destercke. Learning Sets of Probabilities Through Ensemble Methods. 17th European Conference on Symbolic and Quantitative Approaches to Reasoning and Uncertainty (ECSQARU 2023), Sep 2023, Arras, France. pp.270-283, 10.1007/978-3-031-45608-4_21 . hal-04371410

HAL Id: hal-04371410

<https://hal.science/hal-04371410>

Submitted on 3 Jan 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Learning Sets of Probabilities Through Ensemble Methods

Vu-Linh Nguyen^[0000–0003–1642–4468], Haifei Zhang^[0000–0003–4488–1631], and
Sébastien Destercke^[0000–0003–2026–468X]

UMR CNRS 7253, Heudiasyc,
Sorbonne Université, Université de Technologie de Compiègne, France
{vu-linh.nguyen,haifei.zhang,sebastien.destercke}@hds.utc.fr

Abstract. A possible approach to obtain set-valued predictions is to learn for each query instance a probability set (a.k.a. credal set) representing its associated uncertainty. Theoretically founded decision rules extending classical expectation and inducing a partial order between predictions can be used to derive set-valued predictions. However, obtaining such a credal set by imprecisating a given learning algorithm is usually computationally challenging, except for simple models such as decision trees or naive Bayes classifiers. In this paper, we propose a simple, easy to use quantile-based framework for estimating credal sets using output of ensemble methods, that can also cope with complex types of data, such as images and mixed/multimodal data, etc. Experiments are conducted to highlight the usefulness of the proposed framework.

Keywords: Ensemble learning · Credal sets approximation · Set-valued prediction · Quantile-based approach.

1 Introduction

Classification algorithms are usually designed to produce, for each instance, a prediction in the form of a unique element of the set of possible outputs. Under the presence of uncertainty, which is often a consequence of model inadequacy and/or data imperfections (in terms of quality and/or quantity), the model can however be uncertain about its predictions and make unreliable precise predictions. In such a case, it might be more desirable to provide imprecise (or indeterminate) set-valued predictions which aims to balance correctness (the true output is an element of the set-valued prediction) and precision (the cardinality of the set-valued prediction) in some appropriate manner [11,24,34,40].

Learning with a reject option is the simplest case of learning set-valued predictions, in which the classifier is allowed to either produce a singleton prediction or refuse to make a prediction for a given query instance. Threshold-based classifiers have been proposed for that purpose, in which a (global/local) threshold will be employed to decide whether a query instance should be rejected or predicted and then a conventional classifier is called only if the instance should be classified [2,5,7,14,16,17]. Threshold-based classifiers have been developed for

multi-class classification (MCC) [11,24], when the classifier is allowed to return top (locally/globally) ranked classes. While such classifiers are intuitive and easy to implement, they often require reliable estimates of the class probabilities to be performant, which is hard to ensure when information is lacking.

By considering more expressive uncertainty representations, imprecise probabilistic classifiers [6,8,22,39] can provide, at least in theory, more reliable outputs. They are developed based on the assumption that uncertainty is described by a (not necessarily convex) set of probabilities, i.e., a *credal set* [21], a description to which can then be applied theoretically justified decision rules [19,34] to produce set-valued predictions. Moving from a single distribution to a *credal set* is a natural way to model the lack of information, an aspect that unique probabilities can hardly capture. Unfortunately, imprecise probabilistic classifiers often suffer from the limited use to certain types of (tabular) data, as well as from the high computational cost that represent a credal extension of a given learning method. A solution might be to consider the credal set as a neighbourhood of the initial estimated distribution [23,31], yet ensuring the quality of the initial estimated distribution is a challenge itself.

In this paper, we propose a quantile-based framework for estimating credal sets from the output of ensembles [12]. We specifically seek a correctness-precision trade-off when constructing estimates of credal sets, i.e., the estimates are expected to be informative and at the same time not very large. This shall be done by defining “median” of set of distributions and use the “median” to filter out a proportion of “extreme” distributions before forming credal sets. Moreover, we only require the availability of an ensemble of probabilistic classifiers. Thus, the base learner (ensemble) can be freely chosen according to our needs. This flexibility of the proposed approach is remarkably different from existing imprecise probabilistic classifiers. Therefore, we hope to broaden the use of generalized decision rules [19,34] to applications with complex types of data, such as mixed data [10], image/video [37,38] and multimodal data [27,35].

We provide in Section 2 a minimal description of MCC with sets of probabilities. Our main contribution which is a quantile-based approach for estimating credal sets is presented in Section 3. The inference problem with sets of probabilities is summarized in Section 4. Section 5 presents some preliminary experiments on tabular data sets to motivate the use of the proposed framework. Section 6 concludes this work and sketches out future work.

2 Preliminary

We shall recall basics of classification with sets of probabilities and notations.

2.1 Probabilistic Classification

Let \mathcal{X} denote an instance space, and let $\mathcal{Y} = \{y^1, \dots, y^K\}$ be a finite set of classes. We assume that an instance $\mathbf{x} \in \mathcal{X}$ is (probabilistically) associated with members of \mathcal{Y} . We denote by $\mathbf{p}(Y|\mathbf{x})$ the conditional distribution of Y given

$\mathbf{X} = \mathbf{x}$. Given training data $\mathcal{D} = \{(\mathbf{x}_n, y_n) | n = 1, \dots, N\}$ drawn independently from $\mathbf{p}(\mathbf{X}, Y)$, the goal in MCC is to learn a classifier \mathbf{h} , which is a mapping $\mathcal{X} \rightarrow \mathcal{Y}$ that assigns to each instance $\mathbf{x} \in \mathcal{X}$ a class $\hat{y} := \mathbf{h}(\mathbf{x}) \in \mathcal{Y}$.

To evaluate the performance of a classifier \mathbf{h} , a loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ is needed, which compares a prediction \hat{y} with a ground-truth y . Each classifier \mathbf{h} is evaluated using its expected loss

$$R(\mathbf{h}) := \mathbf{E}[\ell(Y, \mathbf{h}(\mathbf{X}))] = \int \ell(y, \mathbf{h}(\mathbf{x})) d\mathbf{P}(\mathbf{x}, y),$$

where \mathbf{P} is the joint probability measure on $\mathcal{X} \times \mathcal{Y}$ characterizing the underlying data-generating process. Therefore, the Bayes-optimal classifier is given by

$$\mathbf{h}^* \in \operatorname{argmin}_{\mathbf{h} \in \mathcal{H}} R(\mathbf{h}), \quad (1)$$

where $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ is the hypothesis space. When \mathcal{H} is probabilistic, we can follow maximum likelihood estimation and define the Bayes-optimal classifier as the classifier which optimizes the conditional log likelihood (CLL) function:

$$\hat{\mathbf{h}} := \hat{\mathbf{p}} \in \operatorname{argmax}_{\mathbf{p} \in \mathcal{H}} \text{CLL}(\mathbf{p} | \mathcal{D}) := \operatorname{argmax}_{\mathbf{p} \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \log \mathbf{p}(y_n | \mathbf{x}_n). \quad (2)$$

To avoid overfitting, the CLL is often augmented by a regularization term [25,29].

Once the classifier (2) is learned from \mathcal{D} , we can in principle find an optimal prediction of any loss function ℓ at the prediction time [13,24]. More precisely, assume the classifier (2) is made available, and predicts for each query instance \mathbf{x} a probability distribution $\mathbf{p}(\cdot | \mathbf{x})$ on the set of labelings \mathcal{Y} . The Bayes-optimal prediction (BOP) of any ℓ is then given by the expected loss minimizer

$$\hat{y} = \hat{y}(\mathbf{x}) \in \operatorname{argmin}_{\bar{y} \in \mathcal{Y}} \mathbf{E}(\ell(y, \bar{y})) = \operatorname{argmin}_{\bar{y} \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} \ell(y, \bar{y}) \mathbf{p}(y | \mathbf{x}). \quad (3)$$

2.2 Classification with Set of Probabilities

Under this setting, we assume that our uncertainty is described by a (not necessarily convex) set of probabilities $\mathcal{P}(\mathcal{Y} | \mathbf{x})$, i.e., a *credal set* [21]. Clearly, the decision rule (3) is no longer directly applicable. Therefore, it is necessary to use some generalized decision rule such as the ones benefiting from strong theoretical justifications [19,34].

Credal sets can arise in different ways, either as a native result of the learning method [1], as the result of an agnostic (with respect to the missingness process) estimation in presence of imprecise data, or as a neighbourhood taken over an initial estimated distribution $\mathbf{p}(Y | \mathbf{x})$ [23,31]. These approaches seems to introduce some inconvenience. Native credal classifiers can be hard to learn, and are unavailable for complex inputs such as mixed data and images. Approximating $\mathcal{P}(\mathcal{Y} | \mathbf{x})$ as a neighbourhood taken over an initial estimated distribution

$\mathbf{p}(Y | \mathbf{x})$ does not face this inconvenience, but requires that the initial estimated distribution is well-estimated, a hard to ensure quality.

In the next section, we propose a simple, flexible and easy to use quantile-based framework for estimating credal sets using output of ensemble methods [12]. This is especially designed to make use of the current and future development of both probabilistic classification and generalized decision rules in a unified framework to broaden the application of imprecise probability (IP) to real-world applications with complex data types.

3 Credal Sets Approximation

We assume an ensemble $\mathbf{H} := \{\mathbf{h}^m | m \in [M] := \{1, \dots, M\}\}$ of M probabilistic classifiers \mathbf{h}^m , $m \in [M]$ is made available and provides, for each instance \mathbf{x} , a set of M probabilistic predictions

$$\mathbf{H}(\mathbf{x}) := \{\mathbf{h}^m(\mathbf{x}) | m \in [M]\} = \{\mathbf{p}^m := (p_1^m, p_2^m, \dots, p_K^m) | m \in [M]\}. \quad (4)$$

Our goal is to aggregate this set of probabilistic predictions into a credal set $\mathcal{P}(\mathcal{Y} | \mathbf{x})$ in some meaningful way.

3.1 A Quantile-Based Approach

The intention of this approach is to seek a correctness-precision trade-off, i.e., the estimations of $\mathcal{P}(\mathcal{Y} | \mathbf{x})$ are expected to be informative and at the same time not very large. We define the reference point of $\mathbf{H}(\mathbf{x})$ as follows:

$$\mathbf{p}^* = \underset{\mathbf{p}: \sum_{k=1}^K p_k = 1}{\operatorname{argmin}} \sum_{m=1}^M d(\mathbf{p}, \mathbf{p}^m). \quad (5)$$

where d is some distance defined for pairs of probability distributions.

Once the reference point \mathbf{p}^* is made available, it allows us to define a preference order, reflecting how common/weird each distribution in $\mathbf{H}(\mathbf{x})$ is:

$$\mathbf{p} \succ \mathbf{p}' \text{ if } d(\mathbf{p}^*, \mathbf{p}) < d(\mathbf{p}^*, \mathbf{p}'). \quad (6)$$

Such a preference order in turn allows us to “discard” a given percentage of outliers among elements of $\mathbf{H}(\mathbf{x})$.

Let $\alpha \in [0, 1]$ be some threshold. We define $\mathbf{H}_\alpha(\mathbf{x})$ as the set of $(1 - \alpha) * 100$ % of closest distributions in $\mathbf{H}(\mathbf{x})$ with respect to the preference order (6). We approximate the credal set $\mathcal{P}(\mathcal{Y} | \mathbf{x})$ of \mathbf{x} by the convex hull of $\mathbf{H}_\alpha(\mathbf{x})$. Let $\mathbf{H}_\alpha(\mathbf{x}) := \{\mathbf{p}^m | m \in [M_\alpha]\}$. The convex hull is defined as

$$\mathbf{CH}_\alpha(\mathbf{x}) := \left\{ \mathbf{p} := \sum_{m=1}^{M_\alpha} \gamma_m \mathbf{p}^m \mid \gamma_m \geq 0, m \in [M_\alpha], \sum_{m=1}^{M_\alpha} \gamma_m = 1 \right\}. \quad (7)$$

The computational complexity of the problem of determining the reference point (5) can greatly depend on the nature of the distance d . In the next section, we recall commonly used distances. Due to page length limit, we only mention few convex distances and refer to [4,15,20,33] for more distances.

3.2 The Cases of Convex Distances

For completeness, we shall start with few definitions and remarks, which are quite basic and would have appeared in textbooks and papers (see, e.g., [3,9,30]).

Definition 1. A function $f : \mathbb{R}^K \mapsto \mathbb{R}$ is convex if for every $\mathbf{p}, \mathbf{p}' \in \mathbb{R}^K$ and every $\lambda_1, \lambda_2 \in [0, 1]$ such that $\lambda_1 + \lambda_2 = 1$, we have the inequality

$$f(\lambda_1 \mathbf{p} + \lambda_2 \mathbf{p}') \leq \lambda_1 f(\mathbf{p}) + \lambda_2 f(\mathbf{p}'). \quad (8)$$

Remark 1. Let $\mathbf{z} \in \mathbb{R}^K$. Let $\|\cdot\|$ be a norm on \mathbb{R}^K . $f(\mathbf{p}) := \|\mathbf{p} - \mathbf{z}\|$ is convex.

Proof. The convexity of $f(\mathbf{p})$ follows consequently from the triangle inequality of norms:

$$\begin{aligned} f(\lambda_1 \mathbf{p} + \lambda_2 \mathbf{p}') &= \|\lambda_1 \mathbf{p} + \lambda_2 \mathbf{p}' - \mathbf{z}\| = \|\lambda_1 (\mathbf{p} - \mathbf{z}) + \lambda_2 (\mathbf{p}' - \mathbf{z})\| \\ &\leq \|\lambda_1 (\mathbf{p} - \mathbf{z})\| + \|\lambda_2 (\mathbf{p}' - \mathbf{z})\| = \lambda_1 \|\mathbf{p} - \mathbf{z}\| + \lambda_2 \|\mathbf{p}' - \mathbf{z}\| \\ &= \lambda_1 f(\mathbf{p}) + \lambda_2 f(\mathbf{p}'). \end{aligned}$$

□

Remark 2. Conical combinations of convex functions are also convex.

Proof. The proof is trivial. It is enough to multiply the inequalities, one per convex function, by non-negative scalars and sum them up. □

In the following, we show that if $f^m(\mathbf{p}) := d(\mathbf{p}, \mathbf{p}^m)$ is convex, $m \in [M]$, then the problem of finding a reference point (5) of $\mathbf{H}(\mathbf{x})$ can be straightforwardly formulated as a convex optimization problem. This is indeed computationally advantageous because with recent advances, convex programming is nearly as straightforward as linear programming [3,32].

Definition 2. A standard convex optimization problem is of the form

$$\underset{\mathbf{p}}{\text{minimize}} \quad f(\mathbf{p}) \quad \text{subject to} \quad g_i(\mathbf{p}) \leq 0, i \in [I], h_j(\mathbf{p}) = 0, j \in [J] \quad (9)$$

where: $\mathbf{p} \in \mathbb{R}^K$ is the optimization variable; The objective function $f : \mathbb{R}^K \mapsto \mathbb{R}$ is convex; The inequality constraint functions $g_i : \mathbb{R}^K \mapsto \mathbb{R}$, $i \in [I]$ are convex; The equality constraint functions $h_j : \mathbb{R}^K \mapsto \mathbb{R}$, $j \in [J]$, are of the form: $h_j(\mathbf{p}) = \mathbf{a}_j \mathbf{p} - b_j$, where \mathbf{a}_j is a vector and b_j is a scalar.

We can encode the condition that the reference point must be a valid probability distribution by using K inequality constraint functions g_i and 1 equality constraint function h_1 :

$$g_k(\mathbf{p}) := -p_k \leq 0, k \in [K], h_1(\mathbf{p}) := \mathbf{1}_K \mathbf{p} - 1 = 0, \quad (10)$$

where $\mathbf{1}_K = (1, \dots, 1)$. The constraints $p_k \leq 1$, $k \in [K]$, are implicitly enforced by the K constraints g_k (i.e., $p_k \geq 0$, $k \in [K]$) and h_1 (i.e., $\sum_{k=1}^K p_k = 1$, $k \in [K]$). Therefore, we can use any existing package to find \mathbf{p}^* (5).

Using Remark 1–2, we can verify that different distances (See [4,15,20,33] and elsewhere) are convex. Examples are members of the L_p Minkowski family

$$f_p(\mathbf{p}) := L_p(\mathbf{p}, \mathbf{z}) := \sqrt[p]{\sum_{k=1}^K |p_k - z_k|^p}, p \geq 1, \quad (11)$$

and Chebyshev distance

$$f_{\text{cheb}}(\mathbf{p}) := L_\infty(\mathbf{p}, \mathbf{z}) := \max_{k \in [K]} |p_k - z_k|. \quad (12)$$

Moreover, a closer look at Definition 1 is enough to verify the convexity of some other distances (discussed in [4,15,20,33] and elsewhere). Examples are the Squared Euclidean distance (whose square function allows triangle inequality)

$$f_{\text{sqe}}(\mathbf{p}) := d^{\text{sqe}}(\mathbf{p}, \mathbf{z}) := \sum_{k=1}^K (p_k - z_k)^2, \quad (13)$$

and KL divergence (inequality (8) can be verified using the log sum inequality):

$$f_{\text{KL}}(\mathbf{p}) := d_{\text{KL}}(\mathbf{p}, \mathbf{z}) := \sum_{k=1}^K p_k \log(p_k / z_k), \quad (14)$$

To solve the problem (9) efficiently, one should carefully look at the nature of the given convex distance. For example, for any given K , closed-form solution for the f_{sqe} (13) can be derived (See Proposition 1). This is also a special case where the additional constraints (i.e., $\sum_{k=1}^K p_k = 1$ and $p_k \geq 0, k \in [K]$) do not change the minimizer. However, it is not always the case. For example, these additional constraints can change the minimizer of f_1 (11) (See Proposition 2). Also, different distances may seek for the same minimizer. Examples of such distances are Topsør and Jensen-Shannon [4]. Moreover, for some distance, such as Inner Product [4], the problem (9) is reduced to a linear program.

Proposition 1. *The reference point \mathbf{p}^* (5) under Squared Euclidean distance f_{sqe} (13) is uniquely defined as*

$$p_k^* = \frac{1}{M} \sum_{m=1}^M p_k^m, k \in [K]. \quad (15)$$

Proof. The proof is trivial and is given for completeness. For any $k \in [K]$, the partial derivative of

$$f(\mathbf{p}) = \sum_{m=1}^M f_{\text{sqe}}^m(\mathbf{p}) = \sum_{k=1}^K \left(\sum_{m=1}^M (p_k - p_k^m) \right)^2 \quad (16)$$

with respect to the variable p_k is

$$\frac{\partial f}{\partial p_k}(\mathbf{p}) = 2 \sum_{m=1}^M (p_k - p_k^m) = 2 \left(M p_k - \sum_{m=1}^M p_k^m \right). \quad (17)$$

Since $f_{\text{sqe}}(\mathbf{p})$ (13) is strictly convex, its unique minimizer is attained when the partial derivatives are zeros, i.e., \mathbf{p}^* is defined in (15). \mathbf{p}^* is a valid distribution because the set of possible distributions is a convex set. \square

Proposition 2. *Except for $K = 2$, the reference point \mathbf{p}^* (5) under f_1 (11) may not be the minimizer of the relaxed optimization problem*

$$\bar{\mathbf{p}} \in \underset{\mathbf{p}}{\operatorname{argmin}} \sum_{m=1}^M L_1(\mathbf{p}, \mathbf{p}^m) = \underset{\mathbf{p}}{\operatorname{argmin}} \sum_{k=1}^K \left(\sum_{m=1}^M |p_k - p_k^m| \right). \quad (18)$$

Proof. Without enforcing the probability axioms (i.e., $\sum_{k=1}^K p_k = 1$ and $p_k \geq 0$, $k \in [K]$), a minimizer $\bar{\mathbf{p}}$ of the relaxed optimization problem (18) is defined as

$$\bar{p}_k := \operatorname{median}(p_k^1, \dots, p_k^M), k \in [K]. \quad (19)$$

This can be verified by showing that, for any $\mathbf{p} \neq \bar{\mathbf{p}}$, we have

$$f_1(p_k) := \sum_{m=1}^M |p_k - p_k^m| \geq \sum_{m=1}^M |\bar{p}_k - p_k^m| := f_1(\bar{p}_k), k \in [K], \quad (20)$$

which implies the relation $f_1(\mathbf{p}) \geq f_1(\bar{\mathbf{p}})$.

Let L_k be the number of p_k^m which is larger than \bar{p}_k . Let S_k be the number of p_k^m which is smaller than \bar{p}_k . By definition of “median”, we have $L_k = S_k$.

– $p_k > \bar{p}_k$: We have the following relations

$$\begin{aligned} |p_k - p_k^m| &= |\bar{p}_k - p_k^m| + |p_k - \bar{p}_k| \text{ if } p_k^m \leq \bar{p}_k, \\ |p_k - p_k^m| &\geq |\bar{p}_k - p_k^m| - |\bar{p}_k - p_k| \text{ if } p_k^m \geq \bar{p}_k. \end{aligned}$$

Therefore, we have

$$\begin{aligned} f_1(p_k) &= \sum_{m=1}^M |p_k - p_k^m| \\ &\geq \sum_{m=1}^M |\bar{p}_k - p_k^m| + |p_k - \bar{p}_k| S_k - |p_k - \bar{p}_k| L_k \\ &= \sum_{m=1}^M |\bar{p}_k - p_k^m| + |p_k - \bar{p}_k| (S_k - L_k) \\ &= \sum_{m=1}^M |\bar{p}_k - p_k^m| = f_1(\bar{p}_k). \end{aligned}$$

– $p_k < \bar{p}_k$: We have the following relations

$$\begin{aligned} |p_k - p_k^m| &= |\bar{p}_k - p_k^m| + |p_k - \bar{p}_k| \text{ if } p_k^m \geq \bar{p}_k, \\ |p_k - p_k^m| &\geq |\bar{p}_k - p_k^m| - |\bar{p}_k - p_k| \text{ if } p_k^m \leq \bar{p}_k. \end{aligned}$$

Therefore, we have

$$\begin{aligned} f_1(p_k) &= \sum_{m=1}^M |p_k - p_k^m| \\ &\geq \sum_{m=1}^M |\bar{p}_k - p_k^m| + |p_k - \bar{p}_k| L_K - |p_k - \bar{p}_k| S_k \\ &= \sum_{m=1}^M |\bar{p}_k - p_k^m| + |p_k - \bar{p}_k| (L_K - S_k) \\ &= \sum_{m=1}^M |\bar{p}_k - p_k^m| = f_1(p'_k). \end{aligned}$$

For $K > 2$, $\bar{\mathbf{p}}$ may not satisfy the probability axioms (see next Table).

$K = 3$					$K > 3$				
\mathbf{p}^1	0.8	0.1	0.1		\mathbf{p}^1	0.4	0.2	$0.4/(K-3)$	$\dots 0.4/(K-3)$
\mathbf{p}^2	0.2	0.5	0.3		\mathbf{p}^2	0.2	0.7	$0.1/(K-3)$	$\dots 0.1/(K-3)$
\mathbf{p}^3	0.1	0.4	0.5		\mathbf{p}^3	0.1	0.6	$0.3/(K-3)$	$\dots 0.3/(K-3)$
$\bar{\mathbf{p}}$	0.2	0.4	0.3		$\bar{\mathbf{p}}$	0.2	0.6	$0.3/(K-3)$	$\dots 0.3/(K-3)$

When $K = 2$, the probability axioms of $\bar{\mathbf{p}}$ are ensured by the fact that the total rank of each distribution \mathbf{p}^m , $m \in [M]$, on the first and the second classes is always $M+1$ (as the masses should sum up to 1). Thus, $\bar{\mathbf{p}}$ is either one element of $\mathbf{H}(\mathbf{x})$ or the average of two elements of $\mathbf{H}(\mathbf{x})$. Let us illustrate this property using an example where $M = 9$:

		\mathbf{p}^1	\mathbf{p}^2	\mathbf{p}^3	\mathbf{p}^4	\mathbf{p}^5	\mathbf{p}^6	\mathbf{p}^7	\mathbf{p}^8	\mathbf{p}^9
p_1	Value	0.9	0.8	0.7	0.6	0.5	0.4	0.3	0.2	0.1
	Rank	1	2	3	4	5	6	7	8	9
p_2	Value	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
	Rank	9	8	7	6	5	4	3	2	1

In this example, the total rank is 10 and $\bar{\mathbf{p}}$ is \mathbf{p}^5 . □

In the next section, we recall the inference problem with credal sets [19,34].

4 Inference Problem

As said, when our uncertainty is described by a credal set $\mathcal{P}(\mathcal{Y}|\mathbf{x})$, instead of a single probability $\mathbf{p}(\mathcal{Y}|\mathbf{x})$, it is necessary to make predictions using some

theoretically founded decision rule extending classical expectation [19,34]. For any $\mathbf{p} \in \mathcal{P}(\mathcal{Y} | \mathbf{x})$ and any loss function ℓ , we shall denote by

$$\hat{y}_\ell^{\mathbf{p}} \in \operatorname{argmin}_{\bar{y} \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} \ell(y, \bar{y}) \mathbf{p}(y | \mathbf{x}). \quad (21)$$

Definition 3. An optimal set-valued prediction under the E-admissibility rule is

$$\hat{\mathbf{Y}}_{\ell, \mathcal{P}}^E = \{y \in \mathcal{Y} | \exists \mathbf{p} \in \mathcal{P} \text{ s.t. } y = \hat{y}_\ell^{\mathbf{p}}\}. \quad (22)$$

Definition 4. An optimal set-valued prediction $\hat{\mathbf{Y}}_{\ell, \mathcal{P}}^M$ under the Maximality rule is the set of the maximal, non-dominated elements of the partial order $\pi_\ell^{\mathcal{P}}$ such that $\bar{y} \succ_{\ell, \mathcal{P}} \bar{y}'$ if

$$\inf_{\mathbf{p} \in \mathcal{P}} \mathbf{E} \mathbf{p} (\ell(y, \bar{y}') - \ell(y, \bar{y})) > 0. \quad (23)$$

In other words, we have

$$\hat{\mathbf{Y}}_{\ell, \mathcal{P}}^M = \{\bar{y} \in \mathcal{Y} | \nexists \bar{y}' \text{ s.t. } \bar{y}' \succ_{\ell, \mathcal{P}} \bar{y}\}. \quad (24)$$

It is known that the set-valued prediction given by the E-admissibility rule is a subset of the one given by the Maximality rule [34].

In the following, we discuss the computational complexity of the inference problem when ℓ is the 0/1 loss, i.e., $\ell(y, \bar{y}) = \llbracket y \neq \bar{y} \rrbracket$, where $\llbracket A \rrbracket = 1$ if the predicate A is true and equals 0 otherwise

Let us start with the case of Maximality rule. For any $\mathbf{p} \in \mathbf{CH}_\alpha(\mathbf{x})$, we have

$$\mathbf{E} \mathbf{p} (\ell(y, \bar{y}') - \ell(y, \bar{y})) = \mathbf{p}(\bar{y}' | \mathbf{x}) - \mathbf{p}(\bar{y} | \mathbf{x}). \quad (25)$$

Thus, the relation $\bar{y} \succ_{\ell, \mathcal{P}} \bar{y}'$ holds if the maximum of the linear program

$$\underset{\mathbf{p}}{\text{maximize}} \quad f(\mathbf{p}) := \mathbf{p}(\bar{y}' | \mathbf{x}) - \mathbf{p}(\bar{y} | \mathbf{x}) \quad (26)$$

$$\text{subject to} \quad \mathbf{p} - \sum_{m=1}^{M_\alpha} \gamma_m \mathbf{p}^m = 0, \gamma_m \geq 0, \sum_{m=1}^{M_\alpha} \gamma_m = 1, \quad (27)$$

is negative. Note that if $f(\mathbf{p})$ has a maximum value on the feasible region, then it has this value on (at least) one of the extreme points, i.e., elements of $\mathbf{H}_\alpha(\mathbf{x})$ [26][Theorem 3.3]. Thus, a naive algorithmic solution is to compute $f(\mathbf{p})$ for the extreme \mathbf{p} and compare it with 0. This requires time $O(K^2 M_\alpha)$ because in the worst case, one needs to check all the $K(K-1)$ relation $\bar{y} \succ_{\ell, \mathcal{P}} \bar{y}'$, $\bar{y} \neq \bar{y}' \in \mathcal{Y}$.

We now tackle the case of the E-admissibility rule. Reminding that, $\forall y \in \hat{\mathbf{Y}}_{\ell, \mathcal{P}}^E$, there must exist at least one $\mathbf{p} \in \mathbf{CH}_\alpha(\mathbf{x})$ such that $y = \hat{y}_\ell^{\mathbf{p}}$. This is equivalent to having at least one $\mathbf{p} \in \mathbf{CH}_\alpha(\mathbf{x})$ such that $\mathbf{p}(y | \mathbf{x}) \geq \mathbf{p}(y' | \mathbf{x})$, $y' \neq y$. Thus, given any outer approximation $\mathcal{Y}_{\ell, \mathcal{P}}^O$ of $\hat{\mathbf{Y}}_{\ell, \mathcal{P}}^E$ we can follow the suggestion of [19] and formulate the problem of checking whether a given $y \in \mathcal{Y}_{\ell, \mathcal{P}}^O$ satisfies

the relation $y \in \hat{\mathbf{Y}}_{\ell, \mathcal{P}}^E$ as finding the maximum value of a linear program

$$\underset{\mathbf{p}}{\text{maximize}} \quad f(\mathbf{p}) := \mathbf{p}(y | \mathbf{x}) - \mathbf{p}(y' | \mathbf{x}) \quad (28)$$

$$\text{subject to} \quad \mathbf{p} - \sum_{m=1}^{M_\alpha} \gamma_m \mathbf{p}^m = 0, \gamma_m \geq 0, \sum_{m=1}^{M_\alpha} \gamma_m = 1, \quad (29)$$

$$\mathbf{p}(y | \mathbf{x}) - \mathbf{p}(y'' | \mathbf{x}) \geq 0, y'' \in \mathcal{Y} \setminus \{y, y'\}, \quad (30)$$

where $y' \neq y$, and comparing it with 0. Hence, finding $\hat{\mathbf{Y}}_{\ell, \mathcal{P}}^E$ requires solving $|\mathcal{Y}_{\ell, \mathcal{P}}^O|$ linear programs, one per $y \in |\mathcal{Y}_{\ell, \mathcal{P}}^O|$. The naive algorithmic solution, i.e., iterating over all the extreme points, can not be applied here because a class y may be optimal only for probabilities in the interior of $\mathbf{CH}_\alpha(\mathbf{x})$.

5 Experiment

To motivate the potential use of the proposed framework, we perform some experiments on 9 tabular datasets from the UCI repository (cf. the left part of Table 1), following a 10-fold cross-validation procedure.

We employ random forests (RFs) [18] (with default setting of scikit-learn) as the base learner. RFs are compared to an instantiation of our framework, where $\mathbf{H}_\alpha(\mathbf{x})$ is constructed under the f_{sqe} (13) and used to produce $\hat{\mathbf{Y}}_{\ell, \mathcal{P}}^M$. For each train test split, we follow a 10-fold nested cross-validation procedure to choose α optimizing u_{65} . The RF is then retrained using the entire training dataset and the chosen α is used to construct $\mathbf{H}_\alpha(\mathbf{x})$ during the inference phase. The source code has been made public at <https://github.com/Haifei-ZHANG/Probability-Sets-Model>.

Table 1. Statistics of data sets (P is the number of features) and experimental results.

Statistics of data sets					Overall results				Cases of abstention	
					RF	Ours			RF	Ours
#	Name	N	P	K	Acc. (%)	u_{65} (%)	u_{80} (%)	$ \hat{\mathbf{Y}}_{\ell, \mathcal{P}}^M $	Acc. (%)	Corr. (%)
1	ecoli	336	7	8	78.35	77.77	79.38	2.05	69.84	93.59
2	balance scale	625	4	3	80.50	82.17	83.15	2.02	26.75	67.75
3	vehicle	846	18	4	74.46	78.16	82.63	2.04	47.31	90.24
4	vowel	990	10	11	65.35	65.89	68.71	2.05	41.05	71.80
5	wine quality	1599	11	6	57.91	61.67	68.54	2.02	49.69	86.73
6	optdigits	1797	64	10	96.95	97.03	97.19	2.03	50.74	80.19
7	segment	2300	19	7	98.05	98.02	98.22	2.09	50.12	78.93
8	waveform	5000	21	3	85.52	85.81	88.33	2	62.06	99.91
9	letter	20000	16	26	96.57	96.58	96.64	2.03	34.33	81.71

Overall results (accuracy, u_{65} and u_{80} scores [40] and cardinality $|\hat{\mathbf{Y}}_{\ell, \mathcal{P}}^M|$) show that our proposal may provide a promising correctness-precision trade-off, compared to RFs itself. Ideally, a reliable classifier should be more cautious on difficult cases, on which the conventional classifier is likely to fail [28,36]. To verify this ability of our proposal, for each dataset, we report the correctness (i.e., the percentage of times the true class is in $\hat{\mathbf{Y}}_{\ell, \mathcal{P}}^M$), given the prediction was imprecise, versus the accuracy of RF on those instances. The results (in the right part of Table 1) strongly support our proposal.

This also suggests that the use of the E-admissibility rule (listed as future work) may improve the overall results because, under the f_{sqe} , predictions of RFs should belong to $\hat{\mathbf{Y}}_{\ell, \mathcal{P}}^E \subset \hat{\mathbf{Y}}_{\ell, \mathcal{P}}^M$ [34]. More precisely, under the f_{sqe} , our proposal should always gain in the term of correctness¹ and the use of the E-admissibility rule may help to produce smaller (reliable) imprecise predictions.

To gain more insights about the influence of α , we consider u_{65} and u_{80} scores on the test set as functions of the value of α . The results in Figure 1 are indeed in agreement with our expectations. The $\mathcal{P}(\mathcal{Y}|\mathbf{x})$ induced by $\mathbf{H}_\alpha(\mathbf{x})$ with small α may contain extreme/noisy distributions and produce large $\hat{\mathbf{Y}}_{\ell, \mathcal{P}}^M$ (resulting in low u_{65} and u_{80} scores). Moderate α may provide a nice correctness-precision trade-off (reflected via promising u_{65} and u_{80} scores). For large α , $\mathcal{P}(\mathcal{Y}|\mathbf{x})$ is shrunk as (small) neighborhood of the \mathbf{p}^* (5) and our proposal (under f_{sqe}) becomes similar to RFs, which use the \mathbf{p}^* to make predictions. The results also suggest that, in practice, nested cross-validation procedure can help us to find some good value of α (even if the ideal gain provided by the optimal α is small).

6 Conclusion

We propose a simple, easy to use quantile-based framework for estimating credal sets using output of ensemble methods, that can also cope with complex types of data. Preliminary experiments suggest that our proposal may provide a promising correctness-precision trade-off, compared to ensemble methods. To seek for a complete picture on the usefulness of our proposal, we envision the following works: (1) implement our proposal with other distances and the E-admissibility rule and analyze (dis)advantages provided by different combinations of distance and decision rule, (2) include threshold-based classifiers as competitors, and (3) include complex types of data (such as images) into our empirical studies.

Our theoretical results also inform that voting ensembles, such as RFs, use the \mathbf{p}^* under f_{sqe} to make predictions. It would be interesting to investigate whether using the \mathbf{p}^* under other distances to make predictions may bring significant difference, though our primary goal is not to study the problem of how to aggregate the probabilistic predictions provided by ensemble members into the final singleton predictions.

¹ Its predictions are identical to the ones provided by RFs in the cases of singleton/precise predictions. In the cases of imprecise predictions, its predictions cover the predictions provided by RFs.

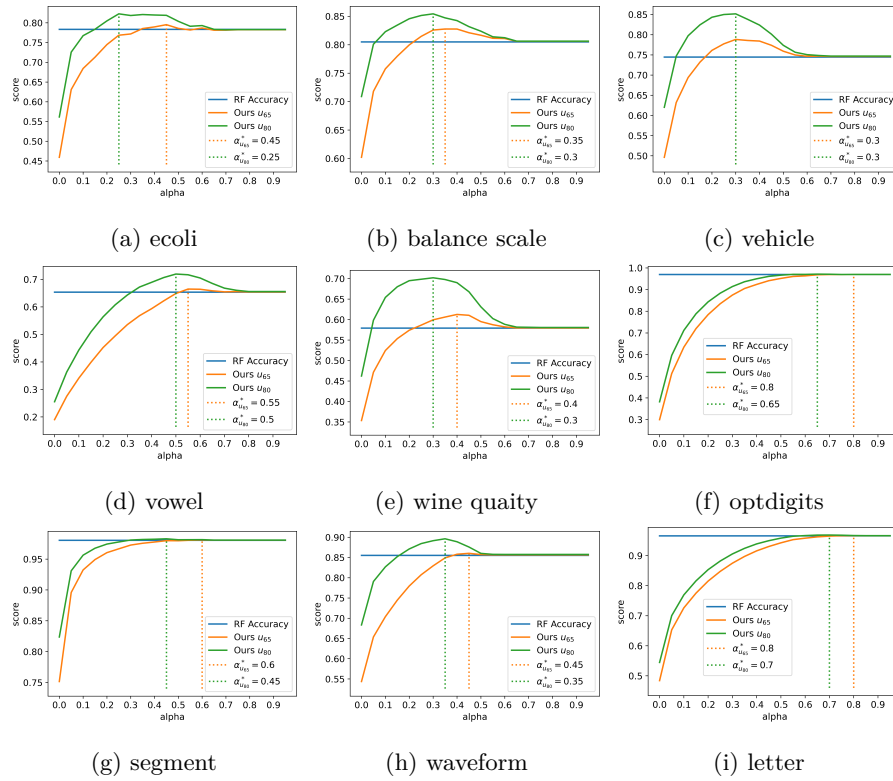


Fig. 1. u_{65} and u_{80} scores on the test set as functions of the value of α

Acknowledgement

This work was funded/supported by the Junior Professor Chair in Trustworthy AI (Ref. ANR-R311CHD).

References

1. Augustin, T., Coolen, F.P., De Cooman, G., Troffaes, M.C.: Introduction to imprecise probabilities. John Wiley & Sons (2014)
2. Bartlett, P.L., Wegkamp, M.H.: Classification with a reject option using a hinge loss. *Journal of Machine Learning Research* **9**(Aug), 1823–1840 (2008)
3. Boyd, S., Boyd, S.P., Vandenberghe, L.: Convex optimization. Cambridge university press (2004)
4. Cha, S.H.: Comprehensive survey on distance/similarity measures between probability density functions. *International Journal of Mathematical models and Methods in Applied Sciences* **1**(4), 300–307 (2007)
5. Chow, C.: On optimum recognition error and reject tradeoff. *IEEE Transactions on information theory* **16**(1), 41–46 (1970)

6. Corani, G., Zaffalon, M.: Learning reliable classifiers from small or incomplete data sets: The naive credal classifier 2. *Journal of Machine Learning Research* **9**(4) (2008)
7. Cortes, C., DeSalvo, G., Mohri, M.: Learning with rejection. In: *Proceedings of the 27th International Conference on Algorithmic Learning Theory (ALT)*. pp. 67–82. Springer Verlag (2016)
8. Cozman, F.G.: Credal networks. *Artificial intelligence* **120**(2), 199–233 (2000)
9. Datta, B.N.: *Numerical linear algebra and applications*, vol. 116. Siam (2010)
10. De Leon, A., Soo, A., Williamson, T.: Classification with discrete and continuous variables via general mixed-data models. *Journal of Applied Statistics* **38**(5), 1021–1032 (2011)
11. Del Coz, J.J., Díez, J., Bahamonde, A.: Learning nondeterministic classifiers. *Journal of Machine Learning Research* **10**(10) (2009)
12. Dietterich, T.G.: Ensemble methods in machine learning. In: *Proceedings of the First International workshop on multiple classifier systems (MCS)*. pp. 1–15. Springer (2000)
13. Elkan, C.: The foundations of cost-sensitive learning. In: *Proceedings of the 17th International Conference on Artificial Intelligence (IJCAI)*. pp. 973–978 (2001)
14. Franc, V., Prusa, D.: On discriminative learning of prediction uncertainty. In: *Proceedings of the 36th International Conference on Machine Learning (ICML)*. pp. 1963–1971 (2019)
15. Gibbs, A.L., Su, F.E.: On choosing and bounding probability metrics. *International statistical review* **70**(3), 419–435 (2002)
16. Grandvalet, Y., Rakotomamonjy, A., Keshet, J., Canu, S.: Support vector machines with a reject option. In: *Proceedings of the 22nd Annual Conference on Neural Information Processing Systems (NIPS)* (2008)
17. Hellman, M.E.: The nearest neighbor classification rule with a reject option. *IEEE Transactions on Systems Science and Cybernetics* **6**(3), 179–185 (1970)
18. Ho, T.K.: Random decision forests. In: *Proceedings of 3rd international conference on document analysis and recognition (ICDAR)*. vol. 1, pp. 278–282. IEEE (1995)
19. Jansen, C., Schollmeyer, G., Augustin, T.: Quantifying degrees of e-admissibility in decision making with imprecise probabilities. In: *Reflections on the Foundations of Probability and Statistics: Essays in Honor of Teddy Seidenfeld*, pp. 319–346. Springer (2022)
20. Lee, L.: Measures of distributional similarity. In: *Proceedings of the 37th annual meeting of the Association for Computational Linguistics on Computational Linguistics (ACL)*. pp. 25–32 (1999)
21. Levi, I.: *The enterprise of knowledge: An essay on knowledge, credal probability, and chance*. MIT press (1983)
22. Mantas, C.J., Abellan, J.: Credal-c4. 5: Decision tree based on imprecise probabilities to classify noisy data. *Expert Systems with Applications* **41**(10), 4625–4637 (2014)
23. Montes, I., Miranda, E., Destercke, S.: Unifying neighbourhood and distortion models: part i—new results on old models. *International Journal of General Systems* **49**(6), 602–635 (2020)
24. Mortier, T., Wydmuch, M., Dembczyński, K., Hüllermeier, E., Waegeman, W.: Efficient set-valued prediction in multi-class classification. *Data Mining and Knowledge Discovery* **35**(4), 1435–1469 (2021)
25. Murphy, K.P.: *Machine learning: a probabilistic perspective*. MIT press (2012)
26. Murty, K.G.: *Linear programming*. Springer (1983)

27. Ngiam, J., Khosla, A., Kim, M., Nam, J., Lee, H., Ng, A.Y.: Multimodal deep learning. In: Proceedings of the 28th international conference on machine learning (ICML). pp. 689–696 (2011)
28. Nguyen, V.L., Destercke, S., Masson, M.H., Hüllermeier, E.: Reliable multi-class classification based on pairwise epistemic and aleatoric uncertainty. In: Proceedings of the 27th International Joint Conference on Artificial Intelligence (IJCAI). pp. 5089–5095 (2018)
29. Nguyen, V.L., Yang, Y., de Campos, C.P.: Probabilistic multi-dimensional classification. In: Proceedings of the 39th Conference on Uncertainty in Artificial Intelligence (UAI). pp. 1522–1533 (2023)
30. Pugh, C.C.: Real mathematical analysis. Undergraduate Texts in Mathematics (2015)
31. Rahimian, H., Mehrotra, S.: Distributionally robust optimization: A review. arXiv preprint arXiv:1908.05659 (2019)
32. Rockafellar, R.T.: Lagrange multipliers and optimality. SIAM review **35**(2), 183–238 (1993)
33. Sriperumbudur, B.K., Gretton, A., Fukumizu, K., Schölkopf, B., Lanckriet, G.R.: Hilbert space embeddings and metrics on probability measures. The Journal of Machine Learning Research **11**, 1517–1561 (2010)
34. Troffaes, M.C.: Decision making under uncertainty using imprecise probabilities. International journal of approximate reasoning **45**(1), 17–29 (2007)
35. Xu, Z., So, D.R., Dai, A.M.: Mufasa: Multimodal fusion architecture search for electronic health records. In: Proceedings of the Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI). vol. 35, pp. 10532–10540 (2021)
36. Yang, G., Destercke, S., Masson, M.H.: Nested dichotomies with probability sets for multi-class classification. In: Proceedings of the Twenty-first European Conference on Artificial Intelligence (ECAI). pp. 363–368 (2014)
37. Yang, Y., Krompass, D., Tresp, V.: Tensor-train recurrent neural networks for video classification. In: Proceedings of the 34th International Conference on Machine Learning (ICML). pp. 3891–3900 (2017)
38. Yin, M., Sui, Y., Liao, S., Yuan, B.: Towards efficient tensor decomposition-based dnn model compression with optimization framework. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 10674–10683 (2021)
39. Zaffalon, M.: The naive credal classifier. Journal of statistical planning and inference **105**(1), 5–21 (2002)
40. Zaffalon, M., Corani, G., Mauá, D.: Evaluating credal classifiers by utility-discounted predictive accuracy. International Journal of Approximate Reasoning **53**(8), 1282–1301 (2012)