



Belief decision trees: theoretical foundations

Zied Elouedi^{a,*}, Khaled Mellouli^a, Philippe Smets^b

^a *Institut Supérieur de Gestion de Tunis, 41 Avenue de la liberté, cité Bouchoucha,
2000 Le Bardo, Tunis, Tunisia*

^b *IRIDIA, Université Libre de Bruxelles, 50 av., F. Roosevelt, CP194/6, 1050 Brussels, Belgium*

Received 1 September 2000; accepted 1 April 2001

Abstract

This paper extends the decision tree technique to an uncertain environment where the uncertainty is represented by belief functions as interpreted in the transferable belief model (TBM). This so-called belief decision tree is a new classification method adapted to uncertain data. We will be concerned with the construction of the belief decision tree from a training set where the knowledge about the instances' classes is represented by belief functions, and its use for the classification of new instances where the knowledge about the attributes' values is represented by belief functions. © 2001 Elsevier Science Inc. All rights reserved.

Keywords: Belief functions; Decision tree; Belief decision tree; Classification; Transferable belief model

1. Introduction

Several learning methods have been developed to ensure classification. Among these, the decision tree method may be one of the most commonly used in supervised learning approaches. Indeed decision trees are characterized by their capability to break down a complex decision problem into several simpler ones. They represent a sequential procedure for deciding the class membership of a given instance. Their major advantage resides in providing a powerful

* Corresponding author.

E-mail addresses: zied.elouedi@isg.rnu.tn (Z. Elouedi), khaled.mellouli@ihc.rnu.tn (K. Mellouli), psmets@ulb.ac.be (P. Smets).

formalism for representing comprehensible classifiers often easy to interpret by experts and even by ordinary users.

Despite their popularity and efficiency, the standard decision trees are unfit to cope with data pervaded with uncertainty both at the construction and at the classification phase. Ignoring the uncertainty may affect classification results and even produce erroneous decisions.

In order to overcome these limitations encountered in standard decision trees, two kinds of techniques have been proposed: *the probabilistic decision trees* [24] and *the fuzzy decision trees* [14,19,22,42,44,45].

The major objective of the probabilistic decision trees is to classify instances with missing or uncertain attribute values where the uncertainty is represented by a probability.

The fuzzy decision trees have been developed to cope with data where both the class and the attributes are represented by fuzzy values. Fuzziness that occurs in classes and attributes is not considered in this paper, but the technique developed here could be extended to such types of data.

The belief function theory, as understood in the transferable belief model (TBM) [36,38], provides a very powerful tool to deal with epistemic uncertainty. It provides a mathematical tool to treat subjective, personal judgments on the different parameters of any classification problem and can be easily extended to deal with objective probabilities. It allows experts to express partial beliefs in a much more flexible way than probability functions do. It permits to handle nicely partial or even total ignorance concerning classification parameters.

Besides, it offers interesting tools to combine several pieces of evidence, like the conjunctive and the disjunctive rules of combination. In addition, decision making is solved through the pignistic transformation.

Hence, the belief function theory seems to provide a convenient framework to handle uncertainty in the decision tree techniques. For this reason, we develop what we call a *belief decision tree approach* which will integrate the advantages of both the decision tree technique and the belief function theory in order to deal with the uncertainty, especially the cognitive one, that may affect the classification problem parameters.

This paper presents basically the theoretical basis of the belief decision trees. It is organized as follows: we start by giving the necessary background concerning the decision trees and the belief function theory. The characteristics of the belief decision trees are then defined in Section 4, whereas in Section 5 the parameters of this new classification method are developed. Finally, in Sections 6 and 7, we detail the construction of a belief decision tree and the classification procedures.

This paper focuses on theoretical foundations of the belief decision trees. Its application and its comparison with other decision trees will be presented in a forthcoming paper. Preliminary results have been published in [8–10]. Another

approach to build belief decision trees is presented in [2,7] where they consider the data as a ‘random sample’ but are limited to binary classes, whereas we use another interpretation of the data as explained in Section 5.

2. Decision trees

A decision tree is a representation of a decision procedure allowing to determine the class of an object. It is composed of three basic elements [26]:

1. *A decision node* specifying the test attribute.
2. *An edge* corresponding to one of the possible values of the test attribute outcomes. It leads generally to a subdecision tree.
3. *A leaf*, which is also named *an answer node*, including objects that, typically, belong to the same class, or at least are very similar.

For what concerns a decision tree, the developer must explain how the tree is constructed and how it is used:

1. *Building the tree.* Based on a given training set, a decision tree is built. It consists in selecting for each decision node the ‘appropriate’ test attribute and also to define the class labeling each leaf.
2. *Classification.* Once the tree is constructed, it is used in order to classify a new instance. We start at the root of the decision tree, we test the attribute specified by this node [23]. The result of this test allows us to move down the tree branch according to the attribute value of the given instance. This process is repeated until a leaf is encountered, the instance then being classified in the same class as the one characterizing the reached leaf.

Several algorithms have been developed in order to ensure the construction of decision trees and its use for the classification task. The ‘ID3’ and ‘C4.5’ algorithms developed by Quinlan [21,25] are probably the most popular ones. We can also mention the ‘CART’ algorithm of Breiman et al. [3].

The majority of the algorithms for building decision trees use a descendent strategy (from the root to the leaves). To ensure this approach, many parameters are required and they can be considered as generic parameters of the algorithm.

The formalism for building decision trees is also referred to as top down induction of decision tree (TDIDT) since it proceeds by successive divisions of the training set where each division represents a question about an attribute value.

A generic decision tree algorithm is characterized by the following properties:

1. *The attribute selection measure.* An attribute is chosen in order to partition the training set in an *optimized* manner. A decision node relative to this attribute is created. It becomes the root of the corresponding decision tree.

2. *A partitioning strategy.* The current training set is divided by taking into account the selected test attributes.
3. *The stopping criteria.* A training subset is declared as a leaf if it satisfies one of the stopping criteria.

The different steps are applied recursively on the training subsets that do not verify the stopping criteria.

Attribute selection measure. This measure allows us to select the attribute that characterizes the root of the decision tree and those of the different sub-decision trees. Quinlan [21] has defined a measure called information gain. It is also referred to as the ‘gain criterion’ based on the information theory of Shannon [30]. The idea is to compute the information gain of each attribute in order to find how well each attribute alone classifies the training examples, then the one presenting the highest value will be chosen. In fact, this attribute generates a partition where the instances classes are as homogeneous as possible within each subset created by the attribute.

Let T denote a training set. Let $\Theta = \{C_1, C_2, \dots, C_n\}$ be the set of n mutually exclusive and exhaustive classes so that each instance in T belongs to one and only one class. Let A be one attribute which domain is finite and denoted by $D(A)$. The *information gain criterion* of Quinlan is defined as the following [25]:

$$\text{Gain}(T, A) = \text{Info}(T) - \text{Info}_A(T), \quad (1)$$

where

$$\text{Info}(T) = - \sum_{i=1}^n \frac{\text{freq}(C_i, T)}{|T|} \cdot \log_2 \frac{\text{freq}(C_i, T)}{|T|}, \quad (2)$$

and

$$\text{Info}_A(T) = \sum_{v \in D(A)} \frac{|T_v|}{|T|} \cdot \text{Info}(T_v), \quad (3)$$

where $\text{freq}(C_i, T)$ denotes the number of objects in the set T that belong to the class C_i and T_v is the subset of objects for which the attribute A has the value v .

The best attribute is the one that maximizes $\text{Gain}(T, A)$. Once the best attribute is allocated to a node, the training set T is split into several subsets, one for each value of the selected attribute. The procedure is then iterated for each subset using only the data that belong to them.

Although it has shown good results, the information gain criterion presents a serious limitation. It favors attributes with a large number of values over those with a small number of values [17,25]. To overcome this drawback, Quinlan has proposed a kind of normalization known as *the gain ratio criterion*. In this manner, the attributes with many values will be adjusted:

$$\text{Gain ratio}(T, A) = \frac{\text{Gain}(T, A)}{\text{Split Info}(T, A)}, \quad (4)$$

where

$$\text{Split Info}(T, A) = - \sum_{v \in D(A)} \frac{|T_v|}{|T|} \cdot \log_2 \frac{|T_v|}{|T|}. \quad (5)$$

$\text{Split Info}(T, A)$ measures the information in the attribute due to the partition of the training set T into $|D(A)|$ training subsets. This quantity describes the information content of the attribute itself.

The idea is to compute the gain ratio of each test attribute, the one presenting the highest value will be selected as the attribute test.

Partitioning strategy. It consists in decomposing the training set into many subsets. In the case of symbolic attributes (with a finite number of values), this strategy resides on testing all the possible attribute values, whereas for the case of numeric attributes, a discretization step is generally needed.

Stopping criteria. It deals with the condition of stopping growth of a part of the decision tree (or even all the decision trees). In other words, it determines whether or not a training subset will be further divided. It is generally fulfilled when all the objects belong to only one class. The part of the decision tree verifying this criterion will be declared as a leaf.

The problems of overfitting and tree pruning are not considered in this paper.

3. Belief function theory

3.1. Basic concepts

In the following, we shall briefly recall some of the basics of the belief function theory. Details can be found in [16,27–29,34,39].

Let Θ be a finite non-empty set including all the elementary events related to a given problem. In the present context, Θ is a set of classes, and the elementary events are the possible classes. These events are assumed to be exhaustive and mutually exclusive. Such a set Θ is classically called the frame of discernment.

The power set of Θ , denoted by 2^Θ , is defined as

$$2^\Theta = \{A : A \subseteq \Theta\}. \quad (6)$$

Hence, the elements of 2^Θ are sets of classes.

The impact of a piece of evidence on the different subsets of the frame of discernment Θ is represented by the so-called basic belief assignment (bba) (initially called basic probability assignment [27], an expression that has

created confusion in the past). The bba m is a function $m : 2^\Theta \rightarrow [0, 1]$ that satisfies

$$\sum_{A \subseteq \Theta} m(A) = 1. \quad (7)$$

The value $m(A)$, called a basic belief mass (bbm), represents the part of belief exactly committed to the subset A of Θ given a piece of evidence, or equivalently to the fact that all we know is that A holds. Due to the lack of information, this quantity cannot be apportioned to any strict subset of A . So, it represents the specific support given to A .

Shafer [27] has initially proposed a normality condition expressed by

$$m(\emptyset) = 0. \quad (8)$$

A bba that satisfies (8) is called a normalized basic belief assignment function.

Smets [33,38] relaxes this condition and interprets $m(\emptyset)$ as the part of belief given to the fact that none of the hypotheses in Θ is true or as the amount of conflict between the pieces of evidence. The subsets A of the frame of discernment Θ , such that $m(A)$ is strictly positive, are called the *focal elements* of the bba m .

A belief function, denoted ‘bel’, corresponding to a specific bba m , assigns to every subset A of Θ the sum of the basic belief masses committed to the subsets of A by m [27]. Contrary to the bba which expresses only the part of beliefs committed exactly to A , the belief function bel represents the total belief that one commits to A without being also committed to \bar{A} . The belief function $\text{bel} : 2^\Theta \rightarrow [0, 1]$ is defined so that

$$\text{bel}(A) = \sum_{\emptyset \neq B \subseteq A} m(B) \quad \forall A \subseteq \Theta. \quad (9)$$

The plausibility function $\text{pl} : 2^\Theta \rightarrow [0, 1]$ quantifies the maximum amount of belief that could be given to a subset A of Θ . It is equal to the sum of the bbm given to subsets B compatible with A . In other words, it contains those parts of beliefs that do not contradict A (i.e., those B such that $B \cap A \neq \emptyset$):

$$\text{pl}(A) = \sum_{A \cap B \neq \emptyset} m(B) \quad \forall A \subseteq \Theta. \quad (10)$$

We also deal with the so-called commonality function $q : 2^\Theta \rightarrow [0, 1]$. It represents the total mass free to move to every element of A [1]. It is defined as

$$q(A) = \sum_{A \subseteq B} m(B) \quad \forall A \subseteq \Theta. \quad (11)$$

Several special belief functions are described:

- The *vacuous belief function* is a belief function that satisfies [27]

$$m(\Theta) = 1 \quad \text{and} \quad m(A) = 0 \quad \forall A \subseteq \Theta, \quad A \neq \Theta. \quad (12)$$

Such bba quantifies the state of *total ignorance* since there is no support given to any strict subset of Θ .

- A *categorical belief function* is a belief function that satisfies [20]

$$m(A) = 1 \text{ for some } A \subseteq \Theta, \quad A \neq \Theta \quad \text{and} \quad m(B) = 0 \quad \forall B \subseteq \Theta, \quad B \neq A. \quad (13)$$

Such a function has a unique focal element A different from the frame of discernment Θ .

- The *certain belief function* is a categorical belief function which focal element is a singleton. It represents a state of total certainty as it assigns all the beliefs to a unique elementary event.

3.2. Combination of belief functions

Let m_1 and m_2 be two bba's defined on the same frame of discernment Θ . These two bba's are induced by two 'distinct' pieces of evidence and collected from two experts (information sources). These bba's can be combined either conjunctively or disjunctively:

1. *The conjunctive rule.* The conjunctive rule of combination handles the case where both sources of information are fully reliable. The result of the combination is a joint bba representing the conjunction of the two pieces of evidence induced from the two sources. Hence, the induced bba quantifies the combined impact of the two pieces of evidence. This rule, denoted by \odot , is defined by [40]

$$(m_1 \odot m_2)(A) = \sum_{B, C \subseteq \Theta: B \cap C = A} m_1(B) \cdot m_2(C). \quad (14)$$

The conjunctive rule can be seen as an unnormalized Dempster's rule of combination. The latter, denoted by \oplus , deals with the closed world assumptions [34], and is defined as [27]

$$(m_1 \oplus m_2)(A) = K \cdot \sum_{B, C \subseteq \Theta: B \cap C = A} m_1(B) \cdot m_2(C), \quad (15)$$

where

$$K^{-1} = 1 - \sum_{B, C \subseteq \Theta: B \cap C = \emptyset} m_1(B) \cdot m_2(C) \quad (16)$$

and

$$(m_1 \oplus m_2)(\emptyset) = 0, \quad (17)$$

K is called *the normalization factor*.

2. *The disjunctive rule.* The dual of the conjunctive rule is the disjunctive rule of combination that builds the bba representing the impact of two pieces of evidence when we only know that at least one is to be accepted, but we do not know which one. This rule, denoted by \oplus , is defined by [40]

$$(m_1 \oplus m_2)(A) = \sum_{B, C \subseteq \Theta: B \cup C = A} m_1(B) \cdot m_2(C). \quad (18)$$

Note that since the conjunctive and the disjunctive rules are both commutative and associative, combining several pieces of evidence induced from distinct information sources (either conjunctively or disjunctively) may be easily ensured by applying repeatedly the chosen rule. The conjunctive and disjunctive rules are not distributive.

3.3. Vacuous extension of belief functions

Let X and Y be two sets of variables such that $Y \subseteq X$. Let m^Y be a bba defined on the domain Θ_Y of Y . The vacuous extension of m^Y to Θ_X , denoted $m^{Y \uparrow X}$, is the bba obtained by extending the information in m^Y to a larger frame X [4,31]:

$$m^{Y \uparrow X}(A \times \Theta_{X-Y}) = m^Y(A) \quad \text{for } A \subseteq \Theta_Y, \quad (19)$$

$$m^{Y \uparrow X}(B) = 0 \quad \text{if } B \text{ is not in the form } A \times \Theta_{X-Y}. \quad (20)$$

3.4. Pignistic transformation

The problem of decision making in the context of the TBM is handled by the pignistic transformation. The TBM is based on a two-level mental model [32]:

- *The credal level* where beliefs are entertained and represented by belief functions.
- *The pignistic level* where beliefs are used to make decisions and represented by probability functions called the pignistic probabilities.

When a decision must be made, beliefs held at the credal level induce a probability measure at the pignistic level, a measure denoted ‘BetP’ [39]. The link between these two functions is achieved by the pignistic transformation:

$$\text{BetP}(A) = \sum_{B \subseteq \Theta} \frac{|A \cap B|}{|B|} \frac{m(B)}{1 - m(\emptyset)} \quad \forall A \subseteq \Theta. \quad (21)$$

4. Characteristics of a belief decision tree

In this section, we introduce the concept of the belief decision tree by presenting its definition and its objectives. Then we define the structure of the

training set which will be illustrated by an example. Finally, the belief decision tree representation will be described.

4.1. Definition

A belief decision tree is a decision tree in an uncertain environment where the uncertainty is represented by belief functions as interpreted in the TBM.

Contrary to a classical decision tree where the objects' classes and the attribute values are known with certainty, in a belief decision tree, these two parameters may be uncertain or even unknown. Such an uncertainty can appear either in the construction or in the classification phase.

4.2. Objectives

As for a standard decision tree, a belief decision tree aims at realizing two major objectives:

1. Building a decision tree from a given set of training instances pervaded with uncertainty. In other words, ensuring the induction of the belief decision tree.
2. Ensuring the classification of new instances which attributes' values can be uncertain or even unknown. Such a procedure is also called the inference procedure.

4.3. Structure of the training set

4.3.1. Definition

Any decision tree is constructed from a training set of objects and based on successive refinements. The training set is the basis leading to the induction of the tree and consequently to the classification of new instances in the inference phase.

This set is traditionally composed of elements represented as pairs (attributes, class) where for each object, we know exactly its assigned class and the value of each attribute.

However, due to the uncertainty introduced here, the structure of the training set we will use may be different from the traditional one. We accept that it may contain data where the class may be uncertain or even unknown. We assume that the values of the attributes of each training instance are known with certainty. Uncertainty on the attributes' values of the cases in the training set is not considered in this paper.

The uncertainty on the classes of the training instance is represented by a basic belief assignment defined on the set of possible classes. This bba, generally given by an expert (or several experts), represents the opinions-beliefs of this expert about the actual value of the class for each object in the training set.

4.3.2. Notations

In this paper, we use the following notations:

T	a given training set
I_j	an instance also named object or case or example
\mathbf{A}	the set of attributes
$\Theta = \{C_1, C_2, \dots, C_n\}$	the frame of discernment involving all the possible classes related to the classification problem; the C_i 's are assumed to be mutually exclusive and exhaustive
$C(I_j)$	the actual class of the object I_j
$m_g^\Theta\{I_j\}[a](C)$	the conditional bbm given to the hypothesis that the actual class $C(I_j)$ of instance I_j belongs to $C \subseteq \Theta$ by an agent g that accepts that the information a is true; useless indices are omitted when they are clearly defined from the context

4.3.3. Special cases

Among the advantages of working under the belief function framework, we notice that the two extreme cases, total knowledge and total ignorance concerning training instance classes, are easily expressed:

- When the class $C(I_j)$ of the object I_j is perfectly known to be $C_v \in \Theta$, it will be represented by a certain belief function having this class as the only focal element:

$$m^\Theta(C_v) = 1 \quad \text{and} \quad m^\Theta(C) = 0 \quad \text{for all } C \neq C_v, \quad C \subseteq \Theta, \quad (22)$$

where $C_v = C(I_j)$ is a singular class.

Such a case is referred to as *total knowledge* and corresponds to the classical *certain* context. Hence, our representation is also appropriate (but not efficient) to describe the standard case where all the classes of the training instances are known with certainty.

- When we have no information about the class of an object which means that the expert is not able to give any judgment concerning the instance's classes, then the bba will be a vacuous belief function defined by

$$m^\Theta(\Theta) = 1 \quad \text{and} \quad m^\Theta(C) = 0 \quad \text{for } C \subset \Theta. \quad (23)$$

Such a case is referred to as *total ignorance*. These cases are generally not included in a training set since they do not provide any information regarding the instances' classes.

- In addition to these special cases, the case of disjunctive classes may be easily described by the so-called categorical belief function. In practice, the latter case may often happen and corresponds to a disjunction of classes that will be assigned to a training instance by an expert. It represents probably the most typical type of information where belief functions will be applied.

Such an opinion will be represented by a bba characterized by only one focal element representing the union of the classes assigned to this training instance.

4.3.4. Example

This is a simple example to illustrate our structure of the training set T within the belief function framework. Let T be a small training set (see Table 1) composed of eight instances. Each training instance is described by three symbolic attributes.

- Hair with possible values {Blond, Red, Dark}.
- Eyes with possible values {Blue, Brown}.
- Height with possible values {Tall, Short}.

The set of the possible classes is $\Theta = \{C_1, C_2, C_3\}$.

So for case 1, the expert strongly supports (0.8) that it is a C_1 case. For case 2, C_2 seems supported but C_1 is not excluded, etc. Methods to assess these bba's are explained in [40].

4.4. Belief decision tree representation

Once the structure of the training set is defined, the representation of our belief decision tree is composed of the same elements as in the traditional decision tree:

- *Decision nodes* for testing attributes.
- *Branches* for specifying attribute values.
- *Leaves* dealing with classes of the training instances.

Due to the uncertainty related to training instances' classes, the structure of leaves will change. Instead of assigning a unique class to each leaf, it will be labeled by a bba expressing a belief on the actual class of the objects belonging to the leaf. The computation of each leaf's bba will be shown in the following section.

Table 1
Training set T^a

Hair	Eyes	Height	bba on classes	C_1	C_2	C_3
Blond	Blue	Tall	$C_1, .8; \Theta, .2$.86	.07	.07
Blond	Brown	Tall	$C_2, .4; C_1 \cup C_2, .4; \Theta, .2$.27	.66	.07
Blond	Blue	Tall	$C_1, .9; \Theta, .1$.94	.03	.03
Blond	Brown	Short	$C_2, .6; C_3, .2; \Theta, .2$.07	.66	.27
Red	Blue	Tall	$C_2, .8; \Theta, .2$.07	.86	.07
Dark	Brown	Short	$C_3, .6; \Theta, .4$.13	.13	.74
Dark	Brown	Tall	$C_3, .9; \Theta, .1$.03	.03	.94
Dark	Brown	Tall	$C_3, .5; C_1 \cup C_3, .2; \Theta, .3$.20	.10	.70

^a Values of the attributes for the eight cases, and the bba's on the classes (represented by pairs which components are the focal element and its bbm). The last three columns give the pignistic probabilities on the three classes.

5. The averaging and the conjunctive approaches for building belief decision trees

The training set in a belief decision tree is characterized by the fact that our knowledge about the value of the actual class of its instances is represented by a bba on Θ , the set of possible classes. The classical algorithms must be adapted to cope with such a context poisoned with uncertainty. Before describing our algorithms in detail, we explain how we derive two of them. Of course, there are many possible algorithms, but interest in the chosen ones comes from their close link to:

- the classical approach developed by Quinlan for our so-called averaging approach;
- the ideas behind the TBM itself for our so-called conjunctive approach.

5.1. The averaging approach

Suppose the decision tree is built, and consider one leaf, denoted S , and suppose that 80% of the cases in S are C_1 cases. Why does this influence our knowledge about the class to which belongs a new case that falls in the same leaf S ? There are (at least) two possible answers.

5.1.1. The sampling answer

We implicitly consider that the new case is a case selected at random from the same population as the one from which the cases in S were selected. We assume that the observed frequencies in S are ‘good estimators’ of the distribution of the classes in the population. Then the proportion 80% is equated (maybe in a not very rigorous way) to the probability that a case randomly selected in the population represented by S is C_1 , i.e., belongs to class C_1 . We observe a new case, and we can say that the probability that the new case is a C_1 is 0.80. Therefore, the 80% proportion of C_1 in S becomes finally the probability that the new case that falls in S is a C_1 .

The largest the dominant proportion in S , the most confident we will feel in our claim about the new case’s class. So we would like that most cases in a leaf belong to the same class. If this ideal is not achievable, we would like to be able to assert with confidence that the new case’s class is one among a few of the possible classes, and that many of the possible classes can be excluded. The worst case for a leaf is encountered when every class is equally represented, as assigning then a class to a new case falling in this leaf would be unjustified and unsupported.

The heterogeneity of the probabilities is what entropy of Eq. (2) is supposed to quantify. The entropy is maximal when the classes are equi-probable, and becomes smaller when the distribution of the classes becomes further and further away from the equi-probability. The smallest entropy (of value 0) is reached when the probability is one for one class, and 0 for all the other classes.

Quinlan's algorithm is based on this idea, and thus tries to minimize the entropy at the leaf's level. Unfortunately, the justification for using the proportions instead of the probabilities can seriously be criticized. It would be acceptable if the number of cases in S was really large, but in practice this is not the case. Leaves with one element are even considered. So the suggested justification is hard to defend, and we feel the second one is more appropriate.

5.1.2. The finite population answer

Suppose that the data of S correspond to a population, not to a sample selected from a larger population. We assume that one case is selected at random, with equi-probability, i.e., with probability $1/|S|$, from S , and that the new case is a duplicate of this selected case, so it is exactly equal to the selected one. If we knew which case had been selected, the class of the new case would be the class of the selected one, but we do not know which one in S was actually selected. All we can say is that the probability that the new case's class is C_i is equal to the probability that the selected case is a C_i , and this probability is equal, thanks to the equi-probable sampling method, to the proportions of C_i cases in S .

Entropy becomes then perfectly meaningful for the same reasons as given in the previous analysis.

For instance, let Table 2 represent the five cases in the subset S of the training set. The class of each case is defined by the indicator function. Cases 1 and 2 are C_1 's, etc. If each case has a probability $p_i = 0.2$ of being selected, then the probability of selecting a C_1 case is just the sum of the indicators weighted by the 0.2 probabilities. In fact the indicator can be understood as the probability that the selected case is a C_j case given the selected case is case i . So the probability that the selected case is C_j becomes

$$P(C_j) = \sum_{i \in S} P(C_j | \text{selected case is } i) p_i. \quad (24)$$

Entropy could then be computed from these 'expected values'.

Table 2
Training subset S (standard case)

Case i	p_i	C_1	C_2	C_3
1	.2	1	0	0
2	.2	1	0	0
3	.2	0	1	0
4	.2	0	1	0
5	.2	0	0	1
Mean		2/5	2/5	1/5

Table 3
Training subset S (probabilistic case)

Case i	p_i	C_1	C_2	C_3
1	.2	.5	.2	.3
2	.2	.4	.6	0
3	.2	0	1	0
4	.2	.1	.7	.2
5	.2	0	.5	.5
Mean		1/5	3/5	1/5

Eq. (24) allows us to shift directly to the case where the classes are uncertain, and *the uncertainty is represented by a probability measure*. Table 3 presents the kind of data that could be collected from the five cases in S . Here case 1 has a high probability of being a C_1 , but might also be a C_2 or a C_3 , even though these last two options are individually less probable than the first. The expected values are computed as in the previous cases, and the probability that the randomly selected case is a C_1 is 0.2 in this case. Entropy could then be computed from these expected values.

What about the entropy in this context? We would like that there would be as few ambiguity as possible when we classify a new case falling in a leaf. So we would like that the probability in a leaf points essentially to one class, and entropy is an excellent measure to quantify this tendency. Hence, the use of the entropy computed from the average probability function in a leaf is plainly justified.

Suppose now that *the uncertainty is represented by a bba*, with m_i the bba about the actual class to which case i belongs. If one were to ask what is the bba of the class to which belongs the randomly selected case, the answer happens to be the average of the m_i 's. This results from the fact the bba's m_i are just conditional bba's and if a case is selected according to a probability distribution, then the resulting bba, denoted \bar{m} , is

$$\bar{m} = \sum_{i \in S} m_i p_i.$$

This formula results from the conjunctive combination of the probability function p_i and the conditional bba's m_i 's, followed by a marginalization of Θ . It uses the relation [35]: $m_1 \odot_2(A) = \sum_{B \subseteq \Omega} m_1[B](A) m_2(B)$, where m_1 and m_2 are defined on Ω . The p_i 's are the m_2 , and the m_i 's are the $m_1[B]$, where $B = \{(i, C_j) : C_j \in \Theta\}$, $\Omega = S \times \Theta$, and $A \subseteq \Theta$.

The probability used to compute the entropy is replaced now by the pignistic probability computed from \bar{m} . It just happens that the pignistic probability computed from \bar{m} is the same as the average of the pignistic probabilities computed for each case:

$$\text{BetP} = \Gamma \left(\sum_{i \in S} m_i p_i \right) = \sum_{i \in S} \Gamma(m_i) p_i,$$

where Γ is the operator that transforms a bba into a pignistic probability function. This linearity property is even the major property that justifies the use of the pignistic probabilities [39]. In consequence, when uncertainty is represented by bba's, it is enough to compute the pignistic probability over Θ , the set of possible classes, for each case, and proceeds as done in Table 3. The use of the entropy at the leaf level is justified just as in the probabilistic case.

Using bba's instead of probabilities on the classes is really not a real issue. One may then raise the question: what is the interest in using the TBM in such a case?

The answer is to be found in dynamic contexts where the beliefs about the classes for each individual can vary with time. New pieces of information could be collected about the data in the training set. In this case the bba's will be adapted by applying the appropriate Dempster's rules, and the change in the pignistic probabilities will not be different from the ones obtained by updating directly the probability model. So using the TBM, even though not essential when the training set is fixed, becomes interesting when our knowledge about the classes of the data in the training set can vary.

Furthermore, let us assume that the value of a needed attribute of a new case is itself uncertain, like the value being v_1 or v_2 . As explained in Section 8, we will compute the bba m_1 from the data in the leaf reached if v_1 was the case, and the bba m_2 from the data in the leaf reached if v_2 was the case. The combination of these two bba's is obtained by the disjunctive rule of combination, something that brings us far away from the pignistic probabilities, and requires the whole TBM apparatus.

So even though in simple cases, the need for the TBM was not essential, it becomes necessary in more complicated contexts.

5.2. The conjunctive approach

The second method we considered is conceptually much closer to the TBM itself. Let us first reconsider what can be done at the leaf's level. In a given leaf S , we suppose we have several cases and ideally they all belong to the same class. So every case in S belongs to the same class, but we do not know which one. Each case provides a bba m_i , $i = 1, \dots, |S|$, that represents what is known from case i about the class to which it belongs, hence to the class that characterizes S . The belief we can build on the class 'common' to those who belong to S is obtained by combining the m_i , for $i \in S$, by the conjunctive combination rule. This idea is based on similar approaches developed by Denoeux [5,41].

So if m_i is the bba of case I_i in the considered leaf S , we compute

$$m_S = \bigoplus_{i \in S} m_i.$$

The bba m_S is what all the cases in leaf S jointly express about the class of those cases that belong to the leaf S .

So the classification of a new case that falls in leaf S is based on this joint bba m_S , and the class is decided using the pignistic probabilities computed from this bba m_S .

Knowing what will be done once the tree is built, let us now shift to its construction. What ‘nice’ property should be satisfied by the cases in a leaf. Ideally, they should belong to the same class. But their actual classes are unknown. Let us consider two cases with the same bba on Θ , the set of possible classes. This is a reasonable property to be satisfied if both cases fall in the same leaf. So we would like all cases in a leaf have bba’s that are ‘close’ to each others. Thus, a distance between bba’s, and in particular between two bba’s, is required.

5.2.1. Distance between bba’s

Let m_1 and m_2 be two bba’s, both defined on Θ . These two bba’s are vectors in a $2^{|\Theta|}$ -dimensional space. A natural distance is the Euclidian distance between the two vectors. Instead of using the bba’s themselves, we could use any vector that is in one-to-one correspondence with the bba, like the bel vector, or the q vector, etc. So let f_i be such a vector where f_i is a function of m_i , which value at $X \subseteq \Theta$ is denoted $f_i(X)$. We are going to show that $f_i(X) = -\ln(q_i(X))$ is an appropriate choice. We define the distance between two instances I_i and I_j belonging to S as

$$d_{ij}^2 = \sum_{X \subseteq \Theta} (f_i(X) - f_j(X))^2.$$

We can then define the distance among the instances within one group S as the average of the distance between pairs of instances in S :

$$D_S^2 = \frac{1}{2s^2} \sum_{i,j \in S} d_{ij}^2.$$

where $s = |S|$.

Ideally this distance should be minimized if the goal is that cases in S are similar to each other. This ‘intra-group’ distance (because computed within one group) has the advantage that minimizing it is equivalent to maximizing the ‘inter-groups’ distances (the one computed among groups), another criterion that could have been advocated.

The intra-group distance can be shown to be equal to

$$\begin{aligned}
 D_S^2 &= \frac{1}{2s^2} \sum_{ij \in S} \sum_X (f_i(X) - f_j(X))^2 \\
 &= \frac{1}{2s^2} \sum_X \sum_{ij \in S} (f_i(X)^2 - 2f_i(X)f_j(X) + f_j(X)^2) \\
 &= \frac{1}{2s^2} \sum_X \left(2s \sum_{i \in S} f_i(X)^2 - 2 \left(\sum_{i \in S} f_i(X) \right)^2 \right) \\
 &= \frac{1}{s} \sum_X \left(\sum_{i \in S} f_i(X)^2 - \frac{1}{s} \left(\sum_{i \in S} f_i(X) \right)^2 \right) \\
 &= \sum_X \frac{1}{s} \sum_{i \in S} (f_i(X) - \overline{f(X)})^2 \\
 &\propto \sum_X \text{variance}(f_i(X)).
 \end{aligned}$$

The distance D_S^2 depends thus on $\overline{f(X)}$. This average \overline{f} can be seen as the function of a bba \overline{m} that has the following property: when ‘added’ s times, the result has the same ‘weight’ as the ‘addition’ of the s bba’s m_i . As far as within one group, all bba’s will be summarized by the result of their combination by the conjunctive rule of combination, the ‘addition’ operation can be seen as applying the conjunctive combination. So we want

$$\bigoplus_{i=1, \dots, s} \overline{m} = \bigoplus_{i=1, \dots, s} m_i.$$

Thanks to the property of the commonality functions, the conjunctive combination rule can be written as a product, and even better, the logarithm of the commonality function of the combination is the sum of the logarithms of the commonality functions entered into the conjunctive combination.

Let κ be defined as minus the logarithm (basis e as the choice of the basis is arbitrary) of q , so let $\kappa_i(X) = -\ln(q_i(X))$ for $X \subseteq \Omega$. Then we define

$$\kappa_{1,2,\dots,s} = \bigoplus_{i=1,\dots,s} \kappa_i.$$

So, we get

$$\kappa_{1,2,\dots,s} = \sum_i \kappa_i.$$

If we take $f(m_i) = \kappa_i$, and $\overline{\kappa} = \frac{1}{s} \kappa_{1,2,\dots,s}$, we have a function f that satisfies the idea that the impact of the s instances in the group under consideration is equal to the impact of s times the ‘average’ case.

Beware that usually $\bar{\kappa}$ is not the commonality function of a belief function, even though $\odot_{i=1,\dots,s}\bar{\kappa}$ is the commonality function of a belief function. In fact, the \odot operator must and can innocuously be extended to ‘generalized belief functions’, i.e., any real function on Θ which coefficients of its Möbius transform (the equivalent of the basic belief masses) add to 1.

So our proposed intra-group distance of instances $i = 1, \dots, s$ becomes

$$D_s^2 = \frac{1}{|S|} \sum_{X \subseteq \Theta} \sum_{i \in S} (\kappa_i(X) - \overline{\kappa(X)})^2, \quad (25)$$

where $\kappa_i(X) = -\ln q\{I_i\}(X)$ and $\overline{\kappa(X)} = \frac{1}{|S|} \sum_{i \in S} \kappa_i(X)$. The case where $q\{I_i\}(\Theta) = 0$ is solved, thanks to the continuity of all involved functions, by putting a very small mass ϵ on Θ , proceeding with the computation and taking the limit for $\epsilon \rightarrow 0$.

6. Belief decision tree parameters

In this section, we define the major parameters leading to the construction of the belief decision tree in both the averaging and the conjunctive approaches. At first, we describe what we have developed as attribute selection measures to ensure the construction of a belief decision tree. Then we present the partitioning strategy and the stopping criteria. Finally, we detail the structure of leaves in belief decision trees.

6.1. Attribute selection measures in a belief decision tree

One of the fundamental parameters in a decision tree (and consequently in a belief decision tree) is *the attribute selection measure*. This measure is used in order to choose *the best* test attribute at each decision node of the tree [11]. It quantifies the class discrimination power of each attribute.

The attribute selection measure has to provide a division of the training set into smaller subsets that are more homogeneous.

The structure of the training set is characterized by data which class is uncertain. This uncertainty is expressed by a bba on the classes domain.

6.1.1. Averaging approach

Under this approach, the attribute selection measure is based on the entropy computed from the average pignistic probability computed from the pignistic probabilities of each instance in the node. We propose the following steps to choose the appropriate attribute:

1. Compute the pignistic probability of each instance I_j in the training set by

$$\text{BetP}^\Theta\{I_j\}(C_i) = \sum_{C_i \in C \subseteq \Theta} \frac{1}{|C|} \frac{m^\Theta\{I_j\}(C)}{1 - m^\Theta\{I_j\}(\emptyset)} \quad \forall C_i \in \Theta. \quad (26)$$

2. Compute the average pignistic probability function $\text{BetP}^\theta\{S\}$ taken over the set of objects S in order to get the average probability on each class:

$$\text{BetP}^\theta\{S\}(C_i) = \frac{1}{|S|} \sum_{I_j \in S} \text{BetP}^\theta\{I_j\}(C_i). \quad (27)$$

3. Compute the entropy of the average pignistic probabilities in S . This value $\text{Info}(S)$ is equal to

$$\text{Info}(S) = - \sum_{i=1}^n \text{BetP}^\theta\{S\}(C_i) \cdot \log_2 \text{BetP}^\theta\{S\}(C_i). \quad (28)$$

4. Select an attribute A . Collect the subset S_v^A made with the cases having v as a value for the attribute A .
5. Compute the average pignistic probability for those cases in subset S_v^A . Let the result be denoted $\text{BetP}^\theta\{S_v^A\}$ for $v \in D(A)$, $A \in \mathbf{A}$.
6. Compute $\text{Info}_A(S)$ using the same definition as suggested by Quinlan, but using the pignistic probabilities instead of the proportions. We get

$$\text{Info}_A(S) = \sum_{v \in D(A)} \frac{|S_v^A|}{|S|} \text{Info}(S_v^A), \quad (29)$$

where $\text{Info}(S_v^A)$ is computed from Eq. (28) using $\text{BetP}^\theta\{S_v^A\}$.

The term $\text{Info}_A(S)$ is equal to the weighed sum of the different $\text{Info}(S_v^A)$ relative to the considered attribute. These $\text{Info}(S_v^A)$ are weighted by the proportion of objects in S_v^A .

7. Compute the information gain provided by the attribute A in the set of objects S such that

$$\text{Gain}(S, A) = \text{Info}(S) - \text{Info}_A(S). \quad (30)$$

8. Using the Split Info, compute the gain ratio relative to the attribute A :

$$\text{Gain Ratio}(S, A) = \frac{\text{Gain}(S, A)}{\text{Split Info}(S, A)}. \quad (31)$$

9. Repeat for every attribute $A \in \mathbf{A}$ and choose the one that maximizes the gain ratio.

Example 1. Let us continue with the previous example (see Table 1). We start by finding $\text{Info}(T)$ relative to this training set. We have to compute the average pignistic probability (see Table 4).

Table 4

Computation of $\text{BetP}^\theta\{T\}$

	C_1	C_2	C_3
$\text{BetP}^\theta\{T\}$	0.32	0.32	0.36

The results induced from the pignistic transformation mean that the average probability that a training instance chosen randomly from T belongs, respectively, to the classes C_1 , C_2 , and C_3 are, respectively, 0.32, 0.32 and 0.36.

These probabilities will be used to compute $\text{Info}(T)$ described as the entropy relative to the whole training set T (see Eq. (28)):

$$\begin{aligned}\text{Info}(T) &= - \sum_{i=1}^3 \text{BetP}(C_i) \cdot \log_2 \text{BetP}(C_i) \\ &= -0.32 * \log_2 0.32 - 0.32 * \log_2 0.32 - 0.36 * \log_2 0.36 = 1.583.\end{aligned}$$

The value 1.583 represents the average amount of information needed to identify the class of an instance in the training set T .

Once the $\text{Info}(T)$ is calculated, we have to look for the $\text{Info}_{\text{Hair}}(T)$, $\text{Info}_{\text{Eyes}}(T)$ and $\text{Info}_{\text{Height}}(T)$. These computations will be ensured by applying Eq. (29).

Let us do the computation for the eye attribute. Let $\text{BetP}^\theta\{T_{\text{Blue}}^{\text{Eyes}}\}$ and $\text{BetP}^\theta\{T_{\text{Brown}}^{\text{Eyes}}\}$ be the average pignistic probability functions relative, respectively, to the objects belonging to T and having, respectively, blue eyes and brown eyes (see Table 5).

The next step consists in the computation of $\text{Info}_{\text{Eyes}}(T)$. By applying Eq. (29), we get

$$\begin{aligned}\text{Info}_{\text{Eyes}}(T) &= -\frac{3}{8} \sum_{i=1}^3 \text{BetP}^\theta\{T_{\text{Blue}}^{\text{Eyes}}\}(C_i) \cdot \log_2 \text{BetP}^\theta\{T_{\text{Blue}}^{\text{Eyes}}\}(C_i) \\ &\quad - \frac{5}{8} \cdot \sum_{i=1}^3 \text{BetP}^\theta\{T_{\text{Brown}}^{\text{Eyes}}\}(C_i) \cdot \log_2 \text{BetP}^\theta\{T_{\text{Brown}}^{\text{Eyes}}\}(C_i) \\ &= 1.299.\end{aligned}$$

Then we compute the information gain (see Eq. (30)), we get, respectively,

$$\text{Gain}(T, \text{Eyes}) = \text{Info}(T) - \text{Info}_{\text{Eyes}}(T) = 1.583 - 1.299 = 0.284,$$

$$\text{Split Info}(T, \text{Eyes}) = 0.95,$$

$$\text{Gain ratio}(T, \text{Eyes}) = 0.3.$$

Table 5

Computation of $\text{BetP}^\theta\{T_{\text{Blue}}^{\text{Eyes}}\}$ and $\text{BetP}^\theta\{T_{\text{Brown}}^{\text{Eyes}}\}$

	C_1	C_2	C_3
$\text{BetP}^\theta\{T_{\text{Blue}}^{\text{Eyes}}\}$	0.62	0.32	0.06
$\text{BetP}^\theta\{T_{\text{Brown}}^{\text{Eyes}}\}$	0.14	0.32	0.54

By applying the same process, we get

$$\text{Gain Ratio}(T, \text{Hair}) = 0.35,$$

$$\text{Gain Ratio}(T, \text{Height}) = 0.11.$$

The attribute that maximizes the gain ratio is the hair attribute. It will be chosen as the root of the decision tree relative to the training set T and branches are created for each of its possible values (Blond, Red, Dark). The same procedure is applied, iteratively, to the cases that fall in each subset that is created according to the values of the selected attribute.

6.1.2. Conjunctive approach

The conjunctive approach uses the intra-group distance D_S^2 (see Eq. (25)) that quantifies for each attribute value how much objects are close to each other.

The selection attribute measure to build a belief decision tree under the conjunctive approach is made of the following steps:

1. For each case in the training set, compute

$$\kappa\{I_j\}(C) = -\ln q^\Theta\{I_j\}(C) \quad \forall C \subseteq \Theta \quad (32)$$

from the bba $m^\Theta\{I_j\}$.

2. For each attribute value v of attribute A , compute the joint $\kappa\{S_v^A\}$ defined on Θ by

$$\kappa\{S_v^A\} = \sum_{I_j \in S_v^A} \kappa\{I_j\}. \quad (33)$$

3. Hence, for each attribute value, the intra-group distance $\text{Sum } D(S_v^A)$ is defined by

$$\text{Sum } D(S_v^A) = \frac{1}{|S_v^A|} \sum_{I_j \in S_v^A} \sum_{X \subseteq \Theta} \left(\kappa\{I_j\}(X) - \frac{1}{|S_v^A|} \kappa\{S_v^A\}(X) \right)^2. \quad (34)$$

4. Once the different $\text{Sum } D(S_v^A)$ are calculated, for each attribute $A \in \mathbf{A}$, compute $\text{Sum } D_A(S)$ representing the weighted sum of the different $\text{Sum } D(S_v^A)$ relative to each value v of the attribute A :

$$\text{Sum } D_A(S) = \sum_{v \in D(A)} \frac{|S_v^A|}{|S|} \text{Sum } D(S_v^A). \quad (35)$$

5. At this level, we may conclude which attribute will be chosen as a root relative to the set of objects S . It consists in selecting the one presenting the minimal $\text{Sum } D_A(S)$. In other words the attribute presenting a partition of objects in which objects are the closest from each other. Nevertheless, we

can proceed in order to take into account the number of possible values for the domain of the attribute.

6. By analogy to our averaging approach, we may also compute $\text{Diff}(S, A)$ defined as the difference between $\text{Sum } D(S)$ and $\text{Sum } D_A(S)$:

$$\text{Diff}(S, A) = \text{Sum } D(S) - \text{Sum } D_A(S), \quad (36)$$

where

$$\text{Sum } D(S) = \frac{1}{|S|} \sum_{I_j \in S} \sum_{X \subseteq \Theta} \left(\kappa\{I_i\}(X) - \frac{1}{|S|} \kappa\{S\}(X) \right)^2. \quad (37)$$

7. Using the Split Info, compute the Diff ratio relative to the attribute A :

$$\text{Diff Ratio}(S, A) = \frac{\text{Diff}(S, A)}{\text{Split Info}(S, A)}. \quad (38)$$

8. Repeat for every attribute $A \in \mathbf{A}$ and choose the one that maximizes the Diff ratio.

Example 2. We use the training set presented in Example 1 and apply the attribute selection measure based on a conjunctive approach in order to find the test attribute.

We start by computing $\text{Sum } D(T)$, the sum of distances separating each training instance to the whole set T . We have

$$\begin{aligned} m\{T\} &= m^\Theta\{I_1\} \odot m^\Theta\{I_2\} \odot \cdots \odot m^\Theta\{I_8\} \\ \text{Sum } D(T) &= \frac{1}{8} \sum_{I_j \in T} D(I_j, T) \\ &= \frac{1}{8} (D(I_1, T) + D(I_2, T) + \cdots + D(I_8, T)) \\ &= 3.14. \end{aligned}$$

We proceed with the computation for the eye attribute. Two values may be possible which are blue eyes or brown eyes. We define the following bba's:

$$m\{T_{\text{Blue}}^{\text{Eyes}}\} = m^\Theta\{I_1\} \odot m^\Theta\{I_3\} \odot m^\Theta\{I_5\}$$

and

$$m\{T_{\text{Brown}}^{\text{Eyes}}\} = m^\Theta\{I_2\} \odot m^\Theta\{I_4\} \odot m^\Theta\{I_6\} \odot m^\Theta\{I_7\} \odot m^\Theta\{I_8\},$$

$$\text{Sum } D_{\text{Eyes}}(T) = \frac{|T_{\text{Blue}}^{\text{Eyes}}|}{|T|} \text{Sum } D(T_{\text{Blue}}^{\text{Eyes}}) + \frac{|T_{\text{Brown}}^{\text{Eyes}}|}{|T|} \text{Sum } D(T_{\text{Brown}}^{\text{Eyes}}) = 2.13.$$

By applying the same process for the other attributes, we get

$$\text{Sum } D_{\text{Hair}}(T) = 1.87,$$

$$\text{Sum } D_{\text{Height}}(T) = 2.79.$$

We have also to compute $\text{Diff}(T, A)$ for $A \in \{\text{Hair}, \text{Eyes}, \text{Height}\}$. We get

$$\text{Diff}(T, \text{Hair}) = 1.27,$$

$$\text{Diff}(T, \text{Eyes}) = 1.01,$$

$$\text{Diff}(T, \text{Height}) = 0.34.$$

Then the computation of the Diff ratio gives as the following results:

$$\text{Diff Ratio}(T, \text{Hair}) = 0.91,$$

$$\text{Diff Ratio}(T, \text{Eyes}) = 1.06,$$

$$\text{Diff Ratio}(T, \text{Height}) = 0.43.$$

The application of the ‘Diff Ratio’ criterion leads to the choice of the eye attribute as an attribute test relative to the training set T .

6.2. Partitioning strategy

The partitioning strategy, also known as the splitting strategy, defines how to split the training set according to the attribute values. Since we deal with symbolic attributes, we create an edge for each value of the attribute chosen as a decision node. Thus, we get several training subsets where each one is relative to one branch and regrouping objects having the same attribute value.

The partitioning strategy for the construction of a belief decision tree is very similar to the partitioning strategy used in the classic tree. This is due to the fact that the uncertainty concerns the classes of the training instances and not the values of their attributes.

6.3. Stopping criteria

The stopping criteria control the process of the construction of the belief decision tree. It allows to stop the development of a path and to declare the node as a leaf. In other words, it determines whether or not a training subset should be further divided.

Three strategies are proposed as a stopping criterion:

1. If the treated node includes only one instance. Hence, the leaf will contain only one object.
2. If the treated node includes instances of which the m_i 's are equal.
3. If there is no further attribute to test. In other words, if all the attributes are split.

4. If the value of the applied attribute selection measure for the remaining attributes is less than or equal to zero which means that the ‘eventual’ partition does not provide a better separation.

In such a case, a leaf will include one or several instances characterized by the same values for the selected attributes but generally having different bba’s on their actual classes.

6.4. Structure of leaves

The leaf in a belief decision tree will be labeled by a basic belief assignment function since the classes of the different training instances are expressed by the means of basic belief assignments.

The major question is how to compute each leaf’s bba? In fact, two cases must be treated: the averaging approach and the conjunctive approach.

6.4.1. Averaging approach for leaves structure determination

Using the averaging approach in the selection attribute measure, the leaf’s bba will be defined as the following:

1. When only one object belongs to the leaf S , the leaf’s bba would be equal to this object’s bba as defined in the training set.
2. When there are many objects attached to the leaf S , the leaf’s bba would be equal to the average of the different basic belief assignment functions relative to these objects:

$$m^\theta\{S\}(C) = \frac{1}{|S|} \sum_{I_j \in S} m^\theta\{I_j\}(C). \quad (39)$$

6.4.2. Conjunctive approach for leaves structure determination

Using the conjunctive approach for the development of the selection attribute measure, the leaf’s bba will be defined as the following:

1. When only one object belongs to the leaf S , the leaf’s bba would be equal to this object’s bba defined on the training set.
2. When there are many objects attached to the leaf S , the leaf’s bba would be equal to the result of the conjunctive combination of these objects’ bba by using the conjunctive rule:

$$m^\theta\{S\} = \bigodot_{I_j \in S} m\{I_j\}. \quad (40)$$

7. Building procedure

The building procedure is also called *the induction task*. It allows to induce a ‘belief’ decision tree in order to use it for the classification task.

Even in an uncertain context, the building of belief decision trees requires a top down approach based on the conquer and divide principle.

7.1. Description

As mentioned, the algorithm to construct a decision tree is based on three major parameters: the attribute selection measure, the partitioning strategy, the stopping criteria. These parameters must take into account the uncertainty encountered in the training set.

In fact, a belief decision tree is constructed from a training set of objects based on successive refinements. These refinements based on both the attribute selection measure and a partitioning strategy lead to small training subsets. The process may be repeated until leaves are encountered. Such nodes have to satisfy the stopping criteria.

7.2. Algorithm of building a belief decision tree

Let T be a training set composed by objects characterized by l symbolic attributes (A_1, A_2, \dots, A_l) and that may belong to the set of classes $\Theta = \{C_1, C_2, \dots, C_n\}$. For each object I_j ($j = 1, \dots, p$) of the training set will correspond a basic belief assignment $m^\Theta\{I_j\}$ expressing the quantity of beliefs exactly committed to the subsets of classes.

Our algorithm, which uses a TDIDT approach, will have the same skeleton as an ID3 algorithm [21]. Besides, our algorithm is considered as generic since it provides two possibilities for selecting the attributes by using either the averaging approach or the conjunctive one.

The different steps of our algorithm of building a belief decision tree are described as follows:

1. Generate the root node of the decision tree including all the objects of the training set T .
2. Choose which approach will be applied: either the averaging approach or the conjunctive one.
3. Verify if this node satisfies or not the stopping criteria (see Section 6.3):
 - If yes, declare it as a leaf node and compute its corresponding bba according to the chosen approach (see Section 6.4).
 - If not, look for the attribute having the highest attribute selection measure (see Section 6.1). This attribute will be designed as the root of the decision tree related to the whole training set.
4. Apply the partitioning strategy (see Section 6.2) by developing an edge for each attribute value chosen as a root. This partition leads to several training subsets.
5. Create a root node relative to each training subset.

6. Repeat the same process for each training subset from the step 3, while verifying the stopping criterion.
7. Stop when all the nodes of the latter level of the tree are leaves.

Example 3. Let us continue with the example proposed in Section 4.3.4. Let us generate the belief decision trees relative, respectively, to the average approach (the one relative to the conjunctive approach may be done in the same manner).

As computed in Example 1, we have found that

$$\text{Gain Ratio}(T, \text{Hair}) = 0.35,$$

$$\text{Gain Ratio}(T, \text{Eyes}) = 0.3,$$

$$\text{Gain Ratio}(T, \text{Height}) = 0.11.$$

Neither of these attributes satisfy the stopping criteria, so we choose the hair attribute as the root of the decision tree relative to the training set T , since it presents the highest gain ratio.

Therefore, branches are created below this root for each of its possible value (Blond, Red, Dark).

We get the belief decision tree shown in Fig. 1.

We notice that the training subset $T_{\text{Red}}^{\text{Eyes}}$ contains only one object. Hence, the stopping criteria are fulfilled for this subset. As a consequence, the node relative to $T_{\text{Red}}^{\text{Eyes}}$ is declared as a leaf and its corresponding bba will be equal to $m^\theta\{I_5\}$.

For the training subsets $T_{\text{Blond}}^{\text{Eyes}}$ and $T_{\text{Dark}}^{\text{Eyes}}$, we apply the same process as we did for the training set T until the stopping criteria hold.

The final belief decision tree induced by our algorithm is given in Fig. 2, where $m^\theta\{I_{13}\}$ and $m^\theta\{I_{78}\}$ are, respectively, the average bba's relative to the objects $\{I_1\}$ and $\{I_3\}$ for $m^\theta\{I_{13}\}$ and to the objects $\{I_7\}$ and $\{I_8\}$ for $m^\theta\{I_{78}\}$.

So we get (see Eq. (39))

$$m^\theta\{I_{13}\}(C_1) = 0.85, \quad m^\theta\{I_1\}(\Theta) = 0.15,$$

$$m^\theta\{I_{78}\}(C_3) = 0.7, \quad m^\theta\{I_{78}\}(C_1 \cup C_3) = 0.1, \quad m^\theta\{I_{78}\}(\Theta) = 0.2.$$

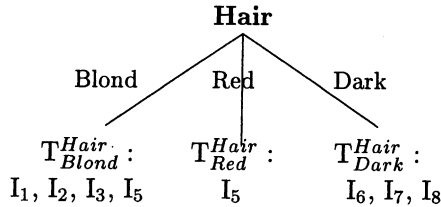


Fig. 1. First generated belief decision tree.

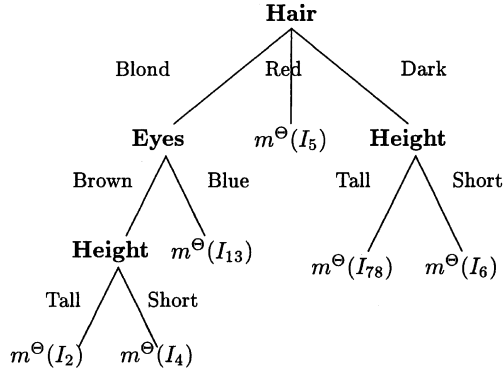


Fig. 2. Belief decision tree representing training instances.

8. Classification procedure

Once the belief decision tree is constructed, the following procedure will be the classification of new instances referring to as new objects. Such a task is also named *the inference task*.

As we deal with an uncertain environment, several cases regarding the knowledge of the attribute values have been studied in order to ensure classification using a belief decision tree.

8.1. Standard classification

Our method is able to ensure the standard classification where each attribute value (of the new instance to classify) is assumed to be exact and certain.

As in an ordinary tree, it consists in starting from the root node and repeating to test the attribute at each node by taking into account the attribute value until reaching a leaf.

Contrary to the classical decision tree where a unique class is attached to the leaf, in our belief decision tree the new instance class will be defined by a basic belief assignment related to the reached leaf. This bba defined on the set of classes represents beliefs on the different subsets of classes (singletons, disjunctions) of the new instance to classify.

In order to make a decision and to get the probability of each singular class, we propose to apply the pignistic transformation.

8.2. Disjunctive case

We consider now the classification of new instances poised by uncertainty in the values of their attributes.

Suppose that the value of an attribute of the new object to classify is not precisely known, and it is only known to belong to a set of possible values of in the domain of this attribute. In particular, when the value is missing, the set is equated to the whole domain of the attribute. This can even occur for several attributes. We assume that the various intervals are non-interactive, i.e., if all we know is that the value of attribute A_1 is in $\theta_1 \subseteq \Theta$ and that the value of attribute A_2 is in $\theta_2 \subseteq \Theta$, then all the values in $\theta_1 \times \theta_2$ are possible. Nevertheless, the algorithm could easily be adapted to the interactive case.

In order to classify such an object, we determine all the leaves which the object could belong to by tracing out all the paths compatible with our knowledge about the different attribute values.

As a consequence, the new instance may belong to many leaves. In each one, we have a bba representing our knowledge about the class to which belong the cases in the leaf. These bba's must be combined in order to get the belief on the instance's class. The disjunctive rule of combination developed by Smets (see Eq. (18)) is the appropriate operator to combine these bba's as it produces the bba under the hypothesis that one path is true (but we do not know which path).

When a decision has to be made, the bba induced from the disjunctive rule is transformed into a probability function by applying the pignistic transformation, producing the probabilities that the new instance belongs to this or that class.

Note that as we mentioned, this case includes the total ignorance of some attribute values. When we deal with unknown attribute values, all the branches relative to the considered attribute will be taken into account. Then the same process (as in the case of the disjunctive values) will be applied.

8.3. General case

The uncertainty characterizing the new instance to classify is not necessarily represented by disjunctive values or missing values (total ignorance). It may be more complicated, especially when the attribute values are provided by several experts. Therefore, it would be interesting to extend our classification procedure in order to handle more general uncertainty.

The uncertainty about the value of the attribute can be defined by a bba on the set of all the possible values of the attribute. Such bba's may be given by one expert or result from the combination (using the conjunctive rule) of several bba's on the attribute values.

Let:

- m^{A_i} be the bba representing the parts of beliefs committed exactly to the different values relative to the attribute A_i of the new instance to classify. This bba is defined on the frame of discernment Θ_{A_i} including all the possible values of the attribute A_i .

- Θ_A be the global frame of discernment relative to all the attributes. It is equal to the cross-product of the different Θ_{A_i} . We denote by

$$\Theta_A = \times_{i=1, \dots, k} \Theta_{A_i}. \quad (41)$$

Since an instance is characterized by a set of combination of values where each one is relative to an attribute, we have firstly to find the bba expressing beliefs on both the different attributes' values of the new instance to classify. In other words, we have to look for the joint bba representing beliefs on all the instance's attributes. To ensure this objective, we have to apply the following steps:

1. Extend the different bba's m^{A_i} to the global frame of attributes Θ_A . As a result, we get the different bba's $m^{A_i \uparrow A}$.
2. Combine the different extended bba's by applying the conjunctive rule of combination:

$$m^{\Theta_A} = \bigodot_{i=1, \dots, k} m^{A_i \uparrow A}. \quad (42)$$

Thus, we get a joint bba representing beliefs on the different combinations of the attributes characterizing the given instance. Should there be some 'correlation' between the bba's, the procedure could be adapted, but in any case the end product is a bba on Θ_A .

We then consider individually the focal elements of the bba m^{Θ_A} . Let x be such a focal element. The next step in our classification task is to compute the belief functions $\text{bel}^\Theta[x]$. The computation of this function depends on the subset x and more exactly on the focal elements of the bba m^{Θ_A} :

1. If the treated focal element x is a singleton (only one value for each attribute), then $\text{bel}^\Theta[x]$ is equal to the belief function corresponding to the leaf to which this focal element is attached.
2. If the focal element x is not a singleton (some attributes have more than one value), i.e., it contains a disjunction in some attribute values. Then we have to explore all the possible paths relative to this combination of values. Two cases are possible:
 - If these paths lead to one leaf, then $\text{bel}^\Theta[x]$ is equal to this leaf's belief function.
 - If these paths lead to distinct leaves, then $\text{bel}^\Theta[x]$ is equal to the result of the disjunctive combination of each leaf's belief function by applying the disjunctive rule.

Finally, the belief functions computed with each focal element x are averaged using the m^{Θ_A} :

$$\text{bel}^\Theta[m^{\Theta_A}](\theta) = \sum_{x \subseteq \Theta_A} m^{\Theta_A}(x) \cdot \text{bel}^\Theta[x](\theta) \quad \text{for } \theta \in \Theta. \quad (43)$$

$\text{bel}^\Theta[m^{\Theta_A}]$ gives us total beliefs on the classes (subsets of classes) that the new instance may to belong. You may also compute the corresponding basic belief assignment $m^\Theta[m^{\Theta_A}]$.

To make a decision, this belief function (or bba) is transformed to a probability function on singular classes via the pignistic transformation.

Example 4. In order to illustrate our generalized case, we consider the instance to classify characterized by:

- Eyes = Blue which is equivalent to a certain bba m^{Eyes} having only the value blue as a focal element: $m^{\text{Eyes}}(\{\text{Blue}\}) = 1$.
- Height = Tall which is equivalent to a certain bba m^{Height} having only the value tall as a focal element: $m^{\text{Height}}(\{\text{Tall}\}) = 1$.
- However, the value relative to the hair attribute is uncertain and described by the bba m^{Hair} such that $m^{\text{Hair}}(\{\text{Dark}\}) = 0.6$; $m^{\text{Hair}}(\{\text{Dark} \cup \text{Red}\}) = 0.2$; $m^{\text{Hair}}(\Theta_{\text{Hair}}) = 0.2$.

Let $\Theta_A = \Theta_{\text{Hair}} \times \Theta_{\text{Eyes}} \times \Theta_{\text{Height}}$.

So $\Theta_A = \{(\text{Blond}, \text{Blue}, \text{Tall}), (\text{Blond}, \text{Blue}, \text{Short}), (\text{Blond}, \text{Brown}, \text{Tall}), (\text{Blond}, \text{Brown}, \text{Short}), (\text{Red}, \text{Blue}, \text{Tall}), (\text{Red}, \text{Blue}, \text{Short}), (\text{Red}, \text{Brown}, \text{Tall}), (\text{Red}, \text{Brown}, \text{Short}), (\text{Dark}, \text{Blue}, \text{Tall}), (\text{Dark}, \text{Blue}, \text{Short}), (\text{Dark}, \text{Brown}, \text{Tall}), (\text{Dark}, \text{Brown}, \text{Short})\}$.

The extension of the different bba's to Θ_A gives as a result:

- The bba m^{Hair} induces the following bba on Θ_A :

$$m^{\text{Hair} \uparrow \Theta_A}(\{\text{Dark}\} \times \Theta_{\text{Eyes}} \times \Theta_{\text{Height}}) = 0.6,$$

$$m^{\text{Hair} \uparrow \Theta_A}(\{\text{Dark} \cup \text{Red}\} \times \Theta_{\text{Eyes}} \times \Theta_{\text{Height}}) = 0.2,$$

$$m^{\text{Hair} \uparrow \Theta_A}(\Theta_{\text{Hair}} \times \Theta_{\text{Eyes}} \times \Theta_{\text{Height}}) = 0.2.$$

- The bba m^{Eyes} induces the following bba on Θ_A :

$$m^{\text{Eyes} \uparrow \Theta_A}(\Theta_{\text{Hair}} \times \{\text{Blue}\} \times \Theta_{\text{Height}}) = 1.$$

- The bba m^{Height} induces the following bba on Θ_A :

$$m^{\text{Height} \uparrow \Theta_A}(\Theta_{\text{Hair}} \times \Theta_{\text{Eyes}} \times \{\text{Tall}\}) = 1.$$

Once the attributes' bba's are extended to Θ_A , then we can apply the conjunctive rule. The result of this combination will be a joint bba on singular instances or subsets of instances. So we get

$$m^{\Theta_A} = m^{\text{Hair} \uparrow \Theta_A} \odot m^{\text{Eyes} \uparrow \Theta_A} \odot m^{\text{Height} \uparrow \Theta_A},$$

where

$$m^{\Theta_A}(\text{Dark} \times \text{Blue} \times \text{Tall}) = 0.6,$$

$$m^{\Theta_A}(\{\text{Dark} \cup \text{Red}\} \times \{\text{Blue}\} \times \{\text{Tall}\}) = 0.2,$$

$$m^{\Theta_A}(\Theta_{\text{Hair}} \times \{\text{Blue}\} \times \{\text{Tall}\}) = 0.2.$$

Next, we have to find beliefs on classes (defined on Θ) given the values of the attributes characterizing the new instance to classify. In fact, three belief functions have to be defined where for each one, we take into account one focal element of m^{Θ_A} . According to the belief decision tree generated (see Fig. 2), we get (see Table 6)

$$\begin{aligned}\text{bel}^{\Theta}[\{(\text{Dark}, \text{Blue}, \text{Tall})\}] &= \text{bel}_{78}, \\ \text{bel}^{\Theta}[\{\text{Dark} \cup \text{Red}\} \times \{\text{Blue}\} \times \{\text{Tall}\}] &= \text{bel}_5 \odot \text{bel}_{78}, \\ \text{bel}^{\Theta}[\Theta_{\text{Hair}} \times \{\text{Blue}\} \times \{\text{Tall}\}] &= \text{bel}_{13} \odot \text{bel}_5 \odot \text{bel}_{78}.\end{aligned}$$

Hence, these belief functions will be averaged (see Eq. (42)), we get:

$$\begin{aligned}\text{bel}^{\Theta}[m^{\Theta_A}](C_1) &= 0.6 * 0 + 0.2 * 0 + 0.2 * 0 = 0, \\ \text{bel}^{\Theta}[m^{\Theta_A}](C_2) &= 0, \\ \text{bel}^{\Theta}[m^{\Theta_A}](C_3) &= 0.6 * 0.7 + 0.2 * 0 + 0 * 0 = 0.42, \\ \text{bel}^{\Theta}[m^{\Theta_A}](C_1 \cup C_2) &= 0, \\ \text{bel}^{\Theta}[m^{\Theta_A}](C_1 \cup C_3) &= 0.6 * 0.8 + 0.2 * 0 + 0.2 * 0 = 0.48, \\ \text{bel}^{\Theta}[m^{\Theta_A}](C_2 \cup C_3) &= 0.6 * 0.7 + 0.2 * 0.56 + 0.2 * 0 = 0.532.\end{aligned}$$

Hence, we get:

$$\begin{aligned}m^{\Theta}[m^{\Theta_A}](C_1) &= 0, \quad m^{\Theta}[m^{\Theta_A}](C_2) = 0, \quad m^{\Theta}[m^{\Theta_A}](C_3) = 0.42, \\ m^{\Theta}[m^{\Theta_A}](C_1 \cup C_2) &= 0, \quad m^{\Theta}[m^{\Theta_A}](C_1 \cup C_3) = 0.06, \\ m^{\Theta}[m^{\Theta_A}](C_2 \cup C_3) &= 0.112, \quad m^{\Theta}[m^{\Theta_A}](\Theta) = 0.408.\end{aligned}$$

Each belief mass represents the part of belief allocated to the fact that the given instance belongs to the focal element of this mass.

Applying the pignistic transformation, the pignistic probability will be defined as follows:

$$\text{BetP}(C_1) = 0.166, \quad \text{BetP}(C_2) = 0.192, \quad \text{BetP}(C_3) = 0.642.$$

We notice that the probability that this instance belongs, respectively, to the classes C_1 , C_2 and C_3 are, respectively, 0.166, 0.192 and 0.642.

Table 6
Beliefs on classes given the attributes' values

	C_1	C_2	C_3	$C_1 \cup C_2$	$C_1 \cup C_3$	$C_2 \cup C_3$	Θ
bel_{78}	0	0	0.7	0	0.8	0.7	1
$\text{bel}_5 \odot \text{bel}_{78}$	0	0	0	0	0	0.56	1
$\text{bel}_{13} \odot \text{bel}_5 \odot \text{bel}_{78}$	0	0	0	0	0	0	1

As a consequence, it is most probable that this new instance (characterized by blue eyes, tall height and uncertain in the color of its hair) belongs to the class C_3 .

9. Conclusion

In this paper, we have defined the belief decision tree which is a new technique associating the decision tree method to the belief function theory in order to handle uncertainty that can exist on classification problem parameters.

We consider the case where the knowledge about the class of the instances in the training set is represented by a belief function over the set of possible classes.

We present two attribute selection measures using the belief function formalism, one parallel to Quinlan's measure based on the entropy (the averaging method), the other close in spirit to the TBM (the conjunctive method). Partitioning strategy and stopping criteria are provided, and the meaning of the data in the leaves is detailed. We present the different steps of the procedure allowing the construction of the tree in an uncertain context.

In fact, these two approaches present two different manners to build belief decision trees. At this level of research, and taking into account theoretical considerations, some differences in the application of these two approaches should be noted.

For the two approaches, the computational complexity is almost the same, but we have to mention that with a very large number of classes, it might become more complex for the conjunctive approach. The two building algorithms are non-incremental and the C4.5 algorithm of Quinlan [25] is a particular case of the averaging approach. Taking into account the underlying hypotheses, the conjunctive approach seems more appropriate since it assumes a total ignorance on instances contrary to the averaging one where instances are considered as equi-probable. The conjunctive approach might be more sensible than the averaging one, but it has the advantage to indicate explicitly this sensitivity via the measure of conflict contrary to the averaging one where conflict is hidden.

Next, we present the inference task ensuring the classification of new instances using the constructed belief decision tree. We consider the case where the knowledge about the value of some attributes is represented by a bba, and show how to use the belief decision tree. All the leaves compatible with the knowledge are considered and their individual conclusions are combined by the disjunctive rule of combination. Classification is then based on the pignistic probabilities derived from the global bba.

This paper presents the theoretical concepts underlying the belief decision trees. The evaluation of the belief decision tree and the comparison of its results with those obtained by classical methods will be reported in a forthcoming paper.

References

- [1] J.A. Barnett, Combining opinions about the order of rule execution, in: The Ninth National Conference on Artificial Intelligence, AAAI, vol. 91, 1991, pp. 477–481.
- [2] M.S. Bjanger, Induction of decision trees from partially classified data using belief functions, Master, University of Compiègne, 2000.
- [3] L. Breiman, J.H. Friedman, R.A. Olshen, C.J. Stone, Classification and Regression Trees, Wadsworth & Brooks, Monterey, CA, 1984.
- [4] J. Cheng, U.M. Fayyad, K.B. Irani, Z. Qian, Improved decision trees: a generalized version of ID3, in: Proceedings of the fifth International Conference on Machine Learning, June 12–14, 1988, pp. 100–106.
- [5] T. Denoeux, A k-nearest neighbor classification rule based on Dempster–Shafer theory, IEEE Transactions on Systems, Man, and Cybernetics 25 (5) (1995) 804–813.
- [6] T. Denoeux, M.S. Bjanger, Induction of decision trees from partially classified data using belief functions, in: IEEE Int. Conf. on Systems, Man and Cybernetics, 2000.
- [7] Z. Elouedi, K. Mellouli, P. Smets, Decision trees using the belief function theory, in: The Proceedings of the International Conference on Information Processing and Management of Uncertainty in Knowledge-based Systems, IPMU'2000, 2000, pp. 141–148.
- [8] Z. Elouedi, K. Mellouli, P. Smets, Classification with belief decision trees, in: The Proceedings of the Ninth International Conference on Artificial Intelligence: Methodology, Systems, Applications, AIMS'2000, 2000, pp. 80–90.
- [9] Z. Elouedi, K. Mellouli, P. Smets, Induction of belief decision trees: a conjunctive approach, in: The Proceedings of the 10th Conference of the Applied Stochastic Models and Data Analysis, ASMDA2001, 2001, pp. 404–409.
- [10] U.M. Fayyad, K.B. Irani, The attribute selection problem in decision tree generation, in: The Tenth National Conference on Artificial Intelligence, AAAI, vol. 92, 1992, pp. 104–110.
- [11] G.Z. Janikow, Fuzzy decision trees: issues and methods, IEEE Transactions on Systems, Man, and Cybernetics, B 28 (1) (1998) 1–14.
- [12] J. Kholas, P.A. Monney, A Mathematical Theory of Hints. An Approach to Dempster–Shafer Theory of Evidence, Lecture Notes on Economics and Mathematical Systems, vol. 425, Springer, Berlin, 1995.
- [13] W.Z. Liu, A.P. White, The importance of attribute selection measures in decision tree induction, Machine Learning 15 (1994) 24–41.
- [14] C. Marsala, Apprentissage inductif en présence de données imprécises: Construction et utilisation d'arbres de décision flous, Thèse de doctorat de l'Université Paris6, LIP6, 1998.
- [15] K. Mellouli, On the propagation of beliefs in networks using the Dempster–Shafer theory of evidence, Ph.D. Dissertation, School of business, University of Kansas, Lawrence, KS, 1987.
- [16] J.R. Quinlan, Induction of decision trees, Machine Learning 1 (1986) 81–106.
- [17] J.R. Quinlan, Decision trees as probabilistic classifiers, in: Proceedings of the Fourth International Workshop on Machine Learning, June 22–25, 1987, pp. 31–37.
- [18] J.R. Quinlan, Decision trees and decision making, IEEE Transactions on Systems, Man, and Cybernetics 20 (2) (1990) 339–346.
- [19] J.R. Quinlan, Probabilistic decision trees, in: R.S. Michalski, J.G. Carbonell, T.M. Mitchell (Eds.), Machine Learning, vol. 3, Morgan Kaufmann, Los Altos, CA, 1990, pp. 140–152 (Chapter 5).
- [20] J.R. Quinlan, C4.5: Programs for Machine Learning, Morgan Kaufmann, San Mateo, CA, 1993.
- [21] S.R. Safavian, D. Landgrebe, A survey of decision tree classifier methodology, IEEE Transactions on Systems, Man, and Cybernetics 21 (3) (1991) 660–674.
- [22] G. Shafer, A Mathematical Theory of Evidence, Princeton University Press, Princeton, NJ, 1976.

- [23] G. Shafer, Perspectives on the theory and practice of belief functions, *International Journal of Approximate Reasoning* 4 (1990) 323–362.
- [24] G. Shafer, Rejoinders to comments on perspectives on the theory and practice of belief functions, *International Journal of Approximate Reasoning* 6 (1992) 445–480.
- [25] C.E. Shannon, A mathematical theory of communication, *The Bell System Technical Journal* 7 (3) (1948) 379–423.
- [26] P.P. Shenoy, G. Shafer, Axioms for probability and belief-function propagation, in: G. Shafer, J. Pearl (Eds.), *The Morgan Kauffmann series in Representation and Reasoning*, 1990, pp. 574–610.
- [27] P. Smets, Decisions and belief functions, Technical Report No. TR/IRIDIA/90-10, Université Libre de Bruxelles, 1990.
- [28] P. Smets, The combination of evidence in the transferable belief model, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 12 (1990) 321–344.
- [29] P. Smets, in: Ph. Smets, A. Mamdani, D. Dubois, H. Prade (Eds.), *Belief Functions, Non Standard Logics for Automated Reasoning*, Academic Press, London, 1990, pp. 253–286.
- [30] P. Smets, Belief functions: the disjunctive rule of combination and the generalized Bayesian theorem, *International Journal of Approximate Reasoning* 9 (1993) 1–35.
- [31] P. Smets, R. Kennes, The transferable belief model, *Artificial Intelligence* 66 (1994) 191–234.
- [32] P. Smets, Numerical representation of uncertainty, in: D.M. Gabbay, Ph. Smets (Eds.), *Handbook of Defeasible Reasoning and Uncertainty Management Systems*, vol. 3, Kluwer Academic Publishers, Dordrecht, 1998, pp. 265–309.
- [33] P. Smets, The transferable belief model for quantified belief representation, in: D.M. Gabbay, Ph. Smets (Eds.), *Handbook of Defeasible Reasoning and Uncertainty Management Systems*, vol. 1, Kluwer Academic Publishers, Dordrecht, 1998, pp. 267–301.
- [34] P. Smets, The application of the transferable belief model to diagnostic problems, *International Journal of Intelligent Systems* 13 (1998) 127–158.
- [35] P. Smets, Practical uses of belief functions, in: K.B. Laskey, H. Prade (Eds.), *Uncertainty in Artificial Intelligence*, UAI99, vol. 15, 1999, pp. 612–621.
- [36] M. Umamo, H. Okamoto, I.H.H. Tamura, F. Kawachi, S. Umedzu, J. Kinoshita, Fuzzy decision trees by fuzzy ID3 algorithm and its application to diagnosis systems, in: *Proceedings of third IEEE International Conference on Fuzzy Systems, FUZZ-IEEE'94*, June 26–29, 1994, pp. 2113–2118.
- [37] Y. Yuan, M.J. Shaw, Induction of fuzzy decision trees, *Fuzzy Sets and Systems* 69 (1995) 125–139.
- [38] J. Zeidler, M. Schlosser, Continuous valued attributes in fuzzy decision trees, *IPMU* 96 (1996) 395–400.