



Apprentissage et classification monotone

RAPPORT DE STAGE

Laura NGUYEN
Licence 3 d'Informatique
Promo 2017-2018

Encadrant : Christophe MARSALA
Équipe : LFI
Laboratoire : LIP6

4 Juin 2018 — 31 Juillet 2018

1 Introduction

Dans beaucoup de problèmes de classification, les valeurs des attributs et de la classe sont ordinaux. De plus, il peut exister une contrainte de monotonie : la classe d'un objet doit croître/décroître en fonction de la valeur de tout ou partie de ses attributs. A savoir, étant donné deux objets x, x' , si $x \leq x'$ alors $f(x) \leq f(x')$. Les variables dépendantes, $f(x)$ et $f(x')$, sont des fonctions monotones des variables indépendantes, x et x' . On parle alors de problèmes de classification monotone, ou problèmes de classification avec contrainte de monotonie. Cette contrainte indique que les objets ayant de meilleures valeurs d'attributs ne doivent pas être assignés à de moins bonnes valeurs de classe.

L'ajout de cette contrainte de monotonie permet d'introduire des concepts sémantiques tels la préférence, la priorité, l'importance, qui nécessitent une relation d'ordre.

Il existe de nombreux domaines se prêtant à ce type de tâches, tels la prédiction du risque de faillite [9], l'analyse de la satisfaction des clients [10], le diagnostic médical [17]. L'importance de la prise en compte d'une relation graduelle entre les valeurs d'attributs et la classe a été démontrée [19] : les classifieurs auxquels sont imposés la contrainte de monotonie sont au moins aussi performants que leurs homologues classiques, et les experts sont plus enclins à utiliser les règles générées par les modèles monotones.

Afin d'extraire des règles à partir de données monotones, on décide d'utiliser les arbres de décision, dont l'efficacité et l'interprétabilité en classification a été prouvée [23]. Cependant, les algorithmes de construction d'arbres de décision standards (générés par CART [16]) ne produisent pas de classifieurs sensibles à la monotonie, même si la base utilisée est complètement monotone. En revanche, il est montré dans [3] que les classifieurs purement monotones ([2], [1], [7]) sont, en terme de taux de bonne classification, statistiquement indiscernables de leurs homologues non-monotones. Dans le même article, il est expliqué que ce phénomène est dû à la sensibilité de ces classifieurs au bruit non-monotone présent dans les données réelles.

Ce stage a pour but d'étudier la construction et l'évaluation d'arbres de décision prenant en compte une relation graduelle susceptible d'exister entre les valeurs d'attributs et la classe, tout en étant suffisamment robuste au bruit non-monotone. On reprend, en particulier, [18] pour la construction d'arbres de décision monotones paramétrés par une mesure de discrimination d'ordre, i.e. une mesure tenant compte de la monotonie d'un attribut par rapport à la classe. Une étude théorique des propriétés des mesures présentées dans le même article est également effectuée.

Dans un premier temps, on présente quelques méthodes présentées dans la littérature. On effectue ensuite l'étude théorique des mesures de discrimination d'ordre, avant de se pencher sur l'implémentation et l'expérimentation de l'algorithme de construction d'arbres monotones.

2 Etat de l'art

Dans cette partie, on présente les méthodes proposées par Hu et al. [13], Marsala et Petturiti [18], Qian et al. [22] et Pei et Hu [21].

2.1 Notations

On considère un ensemble $\Omega = \{\omega_1, \dots, \omega_n\}$ de n d'éléments définis par un ensemble de m attributs $A = \{a_1, \dots, a_m\}$, où pour tout $j = 1, \dots, m$, a_j est une fonction de Ω vers $X_j = \{x_{j1}, \dots, x_{jt_j}\}$. On note aussi $\lambda : \Omega \rightarrow C$ la fonction d'étiquetage, où $C = \{c_1, \dots, c_k\}$ est un ensemble de classes totalement ordonné.

Pour $\omega_i \in \Omega$, l'ensemble dominant de ω_i généré par a_j est défini ([11], [12]) de la façon suivante :

$$[\omega_i]_{a_j}^{\leq} = \{\omega_h \in \Omega : a_j(\omega_i) \leq a_j(\omega_h)\}$$

De même, l'ensemble dominant de ω_i généré par λ s'écrit :

$$[\omega_i]_{\lambda}^{\leq} = \{\omega_h \in \Omega : \lambda(\omega_i) \leq \lambda(\omega_h)\}$$

Les notions de *dominance rough sets* sont également utilisées. On cherche à approximer l'ensemble c_q^{\geq} , i.e. l'ensemble des éléments de Ω dont la classe est inférieure à c_q . Soit $B \subseteq A$ et c_q une classe. Les approximations inférieure et supérieure de c_q^{\geq} sont définies [13] de la façon suivante :

$$R_B^{\geq}(c_q^{\geq}) = \{\omega \in \Omega : [\omega]_B^{\geq} \subseteq c_q^{\geq}\}$$

$$\overline{R_B^{\geq}}(c_q^{\geq}) = \{\omega \in \Omega : [\omega]_B^{\leq} \cap c_q^{\geq} \neq \emptyset\}$$

On note également Ω_α l'ensemble des éléments considérés à chaque étape de l'induction de l'arbre de décision.

2.2 Arbre de décision basé sur l'entropie d'ordre

Il est montré dans [3] que les classifieurs purement monotones ne procurent pas de meilleurs résultats que leurs homologues classiques lorsque les données contiennent un taux élevé de bruit non-monotone.

Afin de résoudre ce problème, les auteurs de [13] proposent un algorithme robuste pour la classification monotone.

2.2.1 Construction du modèle

Pour cela, ils définissent une mesure de degré de monotonie entre deux attributs, qu'ils nomment information mutuelle d'ordre, ou *rank mutual information* (RMI) :

$$RMI^{\leq}(a_j, a_{j'}) = -\frac{1}{n} \sum_{i=1}^n \log \frac{|[\omega_i]_{a_j}^{\leq} \times [\omega_i]_{a_{j'}}^{\leq}|}{n \times |[\omega_i]_{a_j}^{\leq} \cap [\omega_i]_{a_{j'}}^{\leq}|} \quad (1)$$

RMI, sensible à la monotonie et robuste face aux données bruitées, est utilisé comme mesure de sélection d'attribut dans leur algorithme de construction d'arbres de décision monotones, REMT (*Rank Entropy Based Decision Tree*).

Dans cet article, on considère seulement les arbres binaires dans lesquels chaque noeud est affecté à un seul attribut. Concernant les données, les attributs doivent être numériques et la classe, ordinale.

À chaque étape de l'induction :

- si $|\Omega_\alpha| = 1$ ou tous les éléments de Ω_α partagent la même classe, une feuille est créée.
- Sinon, pour chaque attribut a_j , pour chaque valeur d'attribut x_{j_s} , l'ensemble des éléments est divisé en deux sous-ensembles à partir desquels on calcule $RMI_{x_{j_s}} = RMI(a_j^{x_{j_s}}, \lambda)$. On récupère, pour chaque a_j , le seuil de coupure x_j^* maximisant l'information mutuelle d'ordre : $x_{j_s}^* = \arg \max RMI_{x_{j_s}}$. L'attribut a^* utilisé pour le partitionnement est celui dont le seuil x^* maximise $RMI(a_j^{x_{j_s}}, \lambda)$, pour $j = 1, \dots, m, s = 1, \dots, t_j - 1$.
- Si $RMI(a^*, \lambda) < \epsilon$, une feuille est créée.
- Sinon, un noeud est construit et la procédure est répétée sur les sous-ensembles induits par a^* et x^* .

Concernant le critère d'étiquetage, si tous les exemples de la feuille possèdent la même classe, on lui assigne cette dernière. Sinon, si les exemples proviennent de classes différentes, on assigne la classe médiane. Dans le cas où deux classes possèdent le même nombre d'exemples et la feuille courante provient de la branche gauche de son noeud parent, on lui assigne la pire classe. Sinon, on lui affecte la meilleure.

L'approche de construction gloutonne du classifieur ne permet pas d'obtenir un arbre globalement monotone, même si les données utilisées sont monotones. En revanche, cette méthode garantit une monotonie plus faible : il est montré que si la base de données utilisée lors de la construction de l'arbre est monotone-consistante, alors les règles générées par REMT sont monotones [13].

2.2.2 Évaluation du modèle

Afin d'évaluer la performance de leur algorithme, Hu et al. utilisent l'erreur absolue moyenne, ou *Mean Absolute Error* (MAE) :

$$MAE = \frac{1}{n} \sum_{i=1}^n |f(\omega_i) - \lambda(\omega_i)| \quad (2)$$

où $f(\omega_i)$ est la classe prédite par REMT.

2.2.3 Validation du modèle

L'algorithme est testé sur des bases de données artificielles et réelles.

Les bases de données artificielles générées sont monotones. Elles contiennent 1000 exemples à deux attributs, et un nombre de classes variant de 2 à 30.

REMT est ici comparé à CART, Rank Tree [26], OLM [2] et OSDL [6]. Les expériences montrent que REMT produit le moins d'erreurs parmi ces algorithmes, hormis dans le cas à deux classes. Les auteurs soulignent aussi que REMT est le plus précis des quatre algorithmes, peu importe le nombre d'éléments utilisés lors de la construction de l'arbre.

Sur des bases de données réelles, Hu et al. [13] comparent leur modèle à CART et Rank Tree [26]. Les expériences menées montrent que REMT donne de meilleurs résultats en termes de MAE sur 10 bases sur les 12 collectées. De plus, les auteurs montrent également que plus la taille de la base d'apprentissage diminue, plus l'écart de performance se creuse entre REMT et les autres modèles. Enfin, dans le but de comparer les performances des algorithmes lorsque les données sont monotones, les mêmes bases sont modifiées à l'aide d'un algorithme de monotonization [15]. Les résultats des expérimentations montrent que REMT est plus performant que les autres modèles. De plus, les valeurs de MAE produites par tous les modèles diminuent lorsque les données sont monotisées.

2.3 Arbre de décision paramétré par une mesure de discrimination d'ordre

Dans [18], Marsala et Petturiti étudient l'influence de la mesure de discrimination utilisée lors de la construction du classifieur, en termes de taux de bonne prédiction et de monotonie.

Comme les mesures étudiées dans l'article partagent toutes la même structure fonctionnelle, les auteurs définissent un modèle de construction hiérarchique permettant d'isoler les propriétés d'une mesure de discrimination d'ordre et d'en créer de nouvelles.

Soit H^* une mesure de discrimination d'ordre. Le pouvoir de discrimination de l'attribut a_j envers λ selon H^* peut donc s'écrire de la façon suivante :

$$H^*(\lambda|a_j) = h^*(g^*(f^*(\omega_1)), \dots, g^*(f^*(\omega_n)))$$

où f^* mesure la monotonie locale d'un objet par rapport à la classe, g^* est une transformation strictement décroissante de f^* , et h^* agrège les mesures g^* .

Un algorithme de construction d'arbre de décision nommé RDMT (*Rank Discrimination Measure Tree*) est proposé [13]. Cette méthode est essentiellement basée sur REMT mais utilise, à la place de RMI, une mesure de discrimination passée en paramètre.

Comme précédemment, seuls les arbres binaires sont considérés. Ce classifieur gère à la fois les attributs numériques et ordinaux, mais la classe doit être ordinaire. Les données incomplètes ne sont pas acceptées.

2.3.1 Construction du modèle

Notons H la mesure de discrimination utilisée, ϵ le seuil minimal pour H , n_{min} la taille minimale pour Ω_α , δ la profondeur courante de la branche, et Δ la profondeur maximale que peut avoir une branche de l'arbre.

A chaque étape de l'induction :

- si $|\Omega_\alpha| < n_{min}$ ou tous les éléments de Ω_α partagent la même classe ou $\delta > \Delta$, une feuille est créée.
- Sinon, pour chaque attribut a_j , pour chaque valeur d'attribut x_{js} , l'ensemble des éléments est divisé en deux sous-ensembles à partir desquels on calcule $H(\lambda|a_j^{x_{js}})$. On récupère, pour chaque a_j , le seuil de coupure x_j^* minimisant la valeur de H : $x_j^* = \arg \min H(\lambda|a_j^{x_{js}})$. L'attribut a^* utilisé pour le partitionnement est celui dont le seuil x^* minimise $H(\lambda|a_j^{x_{js}})$, pour $j = 1, \dots, m, s = 1, \dots, t_j - 1$.
- Si $H(\lambda|a^*) < \epsilon$, une feuille est créée.
- Sinon, un noeud est construit et la procédure est répétée sur les sous-ensembles induits par a^* et x^* .

Comme REMT, RDMT ne garantit pas un classifieur globalement monotone, mais un arbre générant des règles monotones.

2.3.2 Évaluation du modèle

Les auteurs évaluent les modèles construits selon trois critères : l'accuracy, la monotonie et la taille des arbres (en termes de nombre de feuilles).

L'accuracy est évalué par :

- Le taux de bonne classification (*Correctly Classified Instances*) :

$$CCI = \frac{1}{n} \sum_{i=1}^n (f(\omega_i) \oplus \lambda(\omega_i)) \quad (3)$$

où, $f(\omega_i)$ est la classe prédite pour ω_i par le classifieur, et $f(\omega_i) \oplus \lambda(\omega_i) = 1$ si $f(\omega_i) = \lambda(\omega_i)$, $f(\omega_i) \oplus \lambda(\omega_i) = 0$ sinon.

- **Le test de Kappa**
- MAE (équation 2)

La monotonie d'un arbre est évaluée par deux mesures :

- L'index I de non-monotonie d'un arbre \mathcal{T} :

$$I(\mathcal{T}) = \frac{\sum_{u=1}^q \sum_{v=1}^q m_{u,v}}{q^2 - q} \quad (4)$$

où q est le nombre de feuilles de Γ et $m_{u,v}$ vaut 1 si les feuilles l_u, l_v ne sont pas monotones, 0 sinon.

- L'index NMI1 des exemples dans chaque base de test \mathcal{D} :

$$NMI1(\mathcal{D}) = \frac{\sum_{i=1}^n \sum_{h=1}^n NMP(\omega_i, \omega_h)}{n^2 - n} \quad (5)$$

où $NMP(\omega_i, \omega_h)$ vaut 1 si la paire ω_i, ω_h n'est pas monotone, 0 sinon.

2.3.3 Validation du modèle

Les expériences menées dans l'article visent à comparer les arbres obtenus par différentes mesures ($H_S^*, H_S, H_G^*, H_G, H_P^*, H_{MID}^{10}, H_{ICT}$).

11 bases artificielles sont générées, dont le taux de NMI augmente de 0 à 10% avec un pas de 1%. Chaque base contient 500 exemples à 5 attributs et 5 classes.

Les tests menés sur ces bases montrent qu'il n'existe aucune différence significative entre les mesures en ce qui concerne les indices d'accuracy. De plus, les mesures de discrimination d'ordre produisent des arbres plus monotones, en particulier lorsque le taux de bruit non-monotone de la base augmente. Enfin, il est également montré que ces mesures de discrimination produisent des arbres contenant plus de feuilles que les mesures classiques.

15 bases de données réelles sont sélectionnées pour évaluer les mesures sur des cas d'applications concrets. Comme précédemment, il n'y a pas de différence d'accuracy significative entre les mesures. Cependant, H_S^* , H_G^* , et H_P^* génèrent des arbres plus monotones et comportant davantage de feuilles.

Les arbres produits par des mesures de discrimination d'ordre ne sont pas destinés à être directement utilisés pour la classification monotone. En revanche, étant donné leur degré de monotonie plus élevé, il est avantageux de les passer en entrée d'algorithmes de post-traitement afin de renforcer la monotonie globale.

2.4 Fusion d'arbres de décision monotones

Afin d'améliorer la capacité de généralisation du modèle de classification, Qian et al. proposent dans [22] une méthode d'ensemble par fusion d'arbres de décision monotones, notée FREMT (*Fusing rank entropy based monotonic decision trees*).

Les auteurs couplent une méthode de réduction d'attributs préservant l'ordre avec REMT, utilisé pour générer les classifieurs de base.

Pour réduire l'ensemble des attributs, ils se basent sur les *dominant rough sets* et mettent à jour leur définition. Ils introduisent un paramètre β qui permet de paramétrer la sensibilité du *rough set* au bruit non-monotone. L'ensemble à approximer est :

$$c_q^{\geq} = \bigcup_{u \leq q} C_u$$

où $C_u \in \Omega/\lambda = \{C_1, C_2, \dots, C_r\}$ tel que, pour tout $q, u \leq r$, si $q \geq u$ alors les éléments de C_q sont préférés à ceux de C_u .

Soit $B \subseteq A$ et c_q une valeur de classe. Les approximations inférieure et supérieure de c_q^{\geq} sont définis de la façon suivante :

$$\underline{R}_B^{\geq}(c_q^{\geq}) = \{\omega \in \Omega : \frac{|[\omega]_B^{\geq} \cap c_q^{\geq}|}{|[\omega]_B^{\geq}|} \geq 1 - \beta\} \quad (6)$$

$$\overline{R}_B^{\geq}(c_q^{\geq}) = \{\omega \in \Omega : \frac{|[\omega]_B^{\geq} \cap c_q^{\geq}|}{|[\omega]_B^{\geq}|} \geq \beta\} \quad (7)$$

où $0 \leq \beta \leq 0.5$

La frontière ascendante de c_q liée à l'ensemble d'attributs B s'écrit :

$$BND_B^{\beta}(c_q^{\geq}) = \overline{R}_B^{\geq}(c_q^{\geq}) - \underline{R}_B^{\geq}(c_q^{\geq}) \quad (8)$$

La dépendance monotone de λ selon B est donc exprimée de la façon suivante :

$$\gamma_B^{\beta}(\lambda) = \frac{|\Omega - \bigcup_{q=1}^k BND_B^{\beta}(c_q^{\geq})|}{|\Omega|} \quad (9)$$

A partir de cette définition, il est possible d'établir un coefficient de pertinence d'un attribut a dans B, relativement à λ . La pertinence interne de a dans B relativement à λ est définie de la manière suivante :

$$Sig_{inner}^{\beta}(a, B, \lambda) = \gamma_B^{\beta}(\lambda) - \gamma_{B-\{a\}}^{\beta}(\lambda) \quad (10)$$

De la même façon, $\forall a \in A - B$, la pertinence externe de a selon B correspond à :

$$Sig_{outer}^{\beta}(a, B, \lambda) = \gamma_{B \cup \{a\}}^{\beta}(\lambda) - \gamma_B^{\beta}(\lambda) \quad (11)$$

À l'aide de ces mesures de pertinence d'attributs, Qian et al. [22] proposent un algorithme de réduction d'attributs, paramétré par β , qui, étant donné Ω , A et λ , génère une réduction (*reduct*) d'attributs monotones. Le principe est le suivant :

- On récupère dans B les attributs noyaux de B en utilisant Sig_{inner}^{β} .
- Pour tout attribut $a_j \in A - B$, on calcule $Sig_{outer}^{\beta}(a_j, B, \lambda)$ (équation 11).

- Si la valeur calculée est supérieure ou égale à toutes celles obtenues précédemment, alors on rajoute a_j dans B .
- on réitère le processus jusqu'à ce que, pour tout $a_i \in A - B$, $Sig_{outer}^\beta(a_i, B, \lambda) = 0$ (i.e. rajouter a_i à B ne modifie pas $\gamma_B^\beta(\lambda)$).

Comme la taille de l'ensemble des attributs monotones sélectionnés par cette méthode est bien moindre que celle de l'ensemble original, l'arbre de décision induit par ces attributs sera moins profond et contiendra moins de noeuds, ce qui implique une meilleure aptitude à généraliser. Chaque arbre produit par cet algorithme sert de classifieur de base pour la méthode d'ensembles.

Afin d'obtenir un classifieur final performant, il est nécessaire de produire des arbres les plus diversifiés possibles. Plusieurs sous-ensembles d'attributs sont donc sélectionnés, et chacun de ces sous-ensembles est une réduction d'attributs qui préserve la monotonie des données.

Les auteurs définissent donc une méthode de recherche de multiples réductions d'attributs monotones, paramétrée par β . Étant donné Ω , A et λ , elle consiste à :

- Trouver les attributs noyaux de A en utilisant Sig_{inner}^β (équation 10).
- Trier par ordre croissant les attributs restants $a_j \in B$ en fonction de $\gamma_B^\beta + \frac{|\Omega/\{a_j\}|}{|\Omega|}$.
- Stocker dans P^* les deux sous-ensembles précédents.
- À partir de P^* , trouver une réduction d'attributs RED_0 avec l'algorithme de réduction d'attributs donné plus haut.
- Pour chacun des attributs a_j se trouvant dans RED_0 mais n'étant pas des attributs noyaux,
 - Enlever a_j de P^* .
 - Calculer une réduction RED_j à partir de P^* grâce à l'algorithme précédent.
 - Si RED_j n'a pas déjà été trouvée précédemment, rajouter RED_j à l'ensemble des réductions obtenues.
 - Remettre a_j dans P^* .

Cette méthode permet de générer un ensemble de réductions d'attributs monotones aussi diversifié que possible. Ces sous-ensembles d'attributs servent de base pour construire des arbres de décision complémentaires, générés par REMT.

Soit $F = \{T_1, T_2, \dots, T_N\}$ une forêt d'arbres de décision générés par $RED = \{RED_1, RED_2, \dots, RED_N\}$. Étant donné Ω , A , λ , $\beta = \{\beta_1, \beta_2, \dots\}$, et ω un objet dont la classe est à déterminer, l'algorithme final de fusion d'arbres de décision monotones, *Fusing rank entropy based monotonic decision trees (FREMT)*, est le suivant :

- Pour $i = 1, \dots, M$, trouver $RED_i = \{RED_0, RED_1, \dots, RED_{N_i}\}$ avec l'algorithme précédent
- Stocker dans RED l'union de tous les ensembles de réductions RED_i générés
- Pour chaque $RED_j \in RED$, générer un arbre T_j avec REMT.
- Construire F la forêt d'arbres $\{T_1, T_2, \dots, T_N\}$.
- Déterminer la classe de ω avec F selon le principe de fusion basé sur la probabilité maximale (*Fusing principle based on maximal probability*).

2.4.1 Évaluation et validation du modèle

Qian et al. [22] évaluent la performance de leur modèle selon le taux de bonne classification (équation 3) et MAE (équation 2).

Les tests sont menés uniquement sur des bases réelles.

Sur les 10 bases récupérées, 9 nécessitent un pré-traitement afin d'être adaptés à l'algorithme de fusion d'arbres (en transformant les problèmes de monotonie descendantes en problèmes de monotonie ascendantes). De plus, afin de comparer les performances des modèles lorsque les données sont monotones, les classes des données sont modifiées à l'aide d'un algorithme de monotonization. Les auteurs observent d'abord les performances des arbres constituant la forêt (construits à partir des réductions d'attributs des données d'origine) puis la performance du modèle final. Ils comparent leurs résultats à ceux obtenus par REMT.

D'après les expériences menées, on observe que les arbres de décision générés par les réductions d'attributs offrent de meilleurs résultats en termes de taux de bonne classification et d'erreur

absolue moyenne dans la plupart des cas. De plus, les valeurs générées sont différentes les unes des autres, ce qui satisfait la contrainte de diversité nécessaire au modèle d'ensembles. Enfin, le classifieur produit par FREMT s'avère être meilleur que celui généré par REMT sur toutes les bases considérées, en termes de taux de bonne classification et d'erreur absolue moyenne. FREMT permet aussi d'améliorer la capacité de généralisation de REMT sur ces tâches de classification.

2.5 Arbres de décision partiellement monotones

Les algorithmes étudiés précédemment font l'hypothèse que tous les attributs sont monotones par rapport à la classe. Néanmoins, la plupart des tâches de classification avec contrainte de monotonie comportent deux sortes d'attributs : ceux dont les valeurs sont linéairement ordonnées selon les valeurs de la classe (les critères), et ceux n'ayant pas de relation monotone avec la classe (mais qui permettent tout de même d'améliorer la performance du classifieur). Dans [21], Pei et Hu proposent donc un algorithme de construction d'arbres partiellement monotones permettant de gérer les deux types d'attributs à la fois.

Pour cela, ils proposent une mesure d'inconsistance d'ordre, *rank inconsistency rate* (RIR), qui permet de distinguer les attributs ordinaires des critères et de déterminer les sens des relations graduelles entre les critères et la classe. Étant donné $a_j \in A$ et λ , on définit :

- URIR (*upward rank inconsistency rate*), le RIR ascendant selon a_j :

$$URIR^{\leq}(a_j, \lambda) = -\frac{1}{|\Omega|} \sum_{i=1}^n \log_2 \frac{|[\omega_i]_{a_j}^{\leq}| \times |[\omega_i]_{\lambda}^{\geq}|}{|\Omega| \times |[\omega_i]_{a_j}^{\leq} \cap [\omega_i]_{\lambda}^{\geq}|}$$

- DRIR (*downward rank inconsistency rate*), le RIR descendant selon a_j :

$$DRIR^{\geq}(a_j, \lambda) = -\frac{1}{|\Omega|} \sum_{i=1}^n \log_2 \frac{|[\omega_i]_{a_j}^{\geq}| \times |[\omega_i]_{\lambda}^{\leq}|}{|\Omega| \times |[\omega_i]_{a_j}^{\geq} \cap [\omega_i]_{\lambda}^{\leq}|}$$

a_j croît avec λ lorsque $URIR^{\leq}(a_j, \lambda) = 0$, et décroît avec λ quand $DRIR^{\geq}(a_j, \lambda) = 0$.

La différence entre $URIR^{\leq}(a_j, \lambda)$ et $DRIR^{\geq}(a_j, \lambda)$ est représentée par :

$$diff_{a_j} = URIR^{\leq}(a_j, \lambda) - DRIR^{\geq}(a_j, \lambda)$$

Pei et Hu établissent un seuil $\delta \in [0, 1]$ permettant de décider s'il existe une relation graduelle entre a_j et λ . Si $diff_{a_j} \in [-\delta, \delta]$, alors a_j n'est pas monotone par rapport à λ . Sinon, si $diff_{a_j} \in [\delta, \infty[$, alors a_j décroît avec λ . Au contraire, si $diff_{a_j} \in]-\infty, -\delta]$, alors a_j croît avec λ .

Ces mesures permettent de définir un algorithme capable de déterminer s'il existe une relation graduelle entre chaque attribut et la classe ainsi que son éventuel sens de variation.

Afin de partitionner au mieux les données à chaque étape de construction de l'arbre, les mesures de sélection d'attributs RMI et MI sont utilisées.

L'information mutuelle MI (*Mutual Information*) correspond au degré de consistance entre un attribut et les valeurs de la classe. L'information mutuelle entre l'attribut a_j et la classe λ est défini par :

$$MI(a_j, \lambda) = -\frac{1}{|\Omega|} \sum_{i=1}^n \log_2 \frac{|[\omega_i]_{a_j}| \times |[\omega_i]_{\lambda}|}{|\Omega| \times |[\omega_i]_{a_j} \cap [\omega_i]_{\lambda}|}$$

De plus, Pei et Hu utilisent le coefficient d'information maximale MIC (*Maximal Information Coefficient*) pour enlever les attributs impertinents. Soient X et Y deux variables aléatoires, $|X|$ et $|Y|$ représentent le nombre d'intervalles sur les axes X et Y . Le nombre total d'intervalles n'excède pas une valeur T spécifiée. MIC est défini par :

$$MIC(X, Y) = \max_{|X||Y| < T} \frac{MI(X, Y)}{\log_2(\min(|X|, |Y|))}$$

À l'aide de ces mesures, les auteurs développent PMDT (*Partially Monotonic Decision Trees*), un algorithme de construction d'arbres partiellement monotones. Comme les autres méthodes, celle-ci ne considère que les arbres binaires où à chaque noeud interne correspond un attribut. Elle ne gère

que les attributs numériques, et les valeurs de classe doivent être ordonnées.

Soit Ω_α l'ensemble des exemples considérés à l'étape actuelle de construction, ϵ le seuil minimal pour MI et RMI, et n_{min} la taille minimale pour Ω_α . À chaque étape de l'induction :

- Les éléments $\omega_i \in \Omega$ sont normalisés
- Les attributs impertinents sont retirés avec *MIC* afin d'obtenir Ω'
- Les attributs monotones et les attributs non-monotones sont différenciés en utilisant *RIR*
- Si tous les éléments de Ω' appartiennent à la même classe c , alors une feuille est créée et on lui assigne la classe c
- Sinon, si $|\Omega_\alpha| < n_{min}$ alors une feuille est créée et on lui assigne la classe majoritaire dans Ω_α
- Sinon
 - si le noeud courant n'est pas plus pur que son noeud parent, alors le seuil générant la meilleure partition selon MI est récupéré
 - sinon, le seuil générant la meilleure partition selon RMI est récupéré
- Si la partition obtenue est vide ou la valeur de MI ou RMI est plus petite que ϵ , alors la construction de l'arbre est arrêtée
- Sinon, la même procédure est répétée sur les sous-ensembles de la partition

L'évaluation des modèles est effectuée avec le taux de bonne classification et MAE.

Pei et Hu comparent PMDT à des méthodes basées sur la structure formelle DRSA (*dominance-based rough set approach*) DIR-DOMLEM, VC-DomLEM ([5], [4], [25]), ainsi que des algorithmes de construction d'arbres de décisions tels REMT, RGMT, OLM et OSDL. Ils récupèrent 12 bases réelles dans lesquelles le sens de variation des critères n'est pas connu, et les pré-treatent : normalisation des attributs, suppression des données incomplètes, et transformation des attributs décroissants avec la classe en attributs croissants.

Les expériences menées montrent que PMDT arrive à gérer à la fois les attributs non monotones et les critères, et garde les attributs les plus pertinents pour la construction de l'arbre. MI et RMI gèrent mieux les critères et attributs non-monotones que DomLEM et VC-DomLEM, et permettent d'améliorer significativement les performances de classification en termes de taux de bonne classification et de MAE. De plus, en considérant le fait que les données puissent être totalement monotones ou non, PMDT produit également de meilleurs résultats que REMT, RGMT, OLM et OSDL sur la plupart des tâches de classification traitées.

3 Etude théorique des propriétés des mesures de discrimination d'ordre

Marsala et Petturiti définissent dans [18] un modèle de construction hiérarchique de mesures de discrimination d'ordre. Ce modèle a pour but d'isoler les propriétés qu'une fonction doit avoir pour être une telle mesure. Il sert aussi de base pour en créer de nouvelles.

Soit H^* une mesure de discrimination d'ordre. On rappelle que le pouvoir de discrimination de l'attribut a_j envers la classe λ peut se décomposer de la façon suivante :

$$H^*(\lambda|a_j) = h^*(g^*(f^*(\omega_1), \dots, g^*(f^*(\omega_n))))$$

où f^* est une mesure de la monotonie locale de l'objet, g^* une mesure de non-monotonie de l'objet, et h^* une agrégation des mesures g^* .

Pour que H^* soit une mesure de discrimination d'ordre, chaque couche doit satisfaire certaines conditions.

Soient $\omega_i \in \Omega$ et f^*, g^*, h^* les trois couches de H^* . Les propriétés suivantes sont démontrées en annexe :

4 Implémentation et expérimentation de l'algorithme de construction d'arbres monotones

On implémente une librairie de fonctions permettant la construction ainsi que l'étude expérimentale d'arbres de décision. On se base essentiellement sur RDMT(H), donné dans [18], pour construire des arbres de décision monotones, que l'on évalue sur des données artificielles et réelles. Le code est implémenté en Python 3.6 et les librairies numpy [24], matplotlib [14] et scikit-learn [20] sont utilisées.

4.1 Mesures de discrimination d'ordre

On implémente le modèle de construction hiérarchique proposé par Marsala et Petturiti. On rappelle que le pouvoir de discrimination de l'attribut a_j envers la classe λ selon une mesure de discrimination H^* peut s'écrire de la façon suivante :

$$H^*(\lambda|a_j) = h^*(g^*(f^*(\omega_1)), \dots, g^*(f^*(\omega_n)))$$

Toutes les fonctions f^* , g^* , h^* présentées dans l'article sont codées. Pour chaque fonction f^* , on distingue le cas où la mesure construite détecte la monotonie (i.e. elle considère les ensembles dominants) du cas contraire (i.e. elle considère les classes d'équivalences).

Par exemple, étant donnés $a_j \in \mathcal{A}$ et λ , l'entropie conditionnelle de Shannon se réécrit de la façon suivante :

$$H_S(\lambda|a_j) = \sum_{i=1}^n \frac{1}{n} (-\log_2(\frac{|[\omega_i]_\lambda \cap [\omega_i]_{a_j}|}{|[\omega_i]_{a_j}|}))$$

$$\text{où } f^*(\omega_i) = \frac{|[\omega_i]_\lambda \cap [\omega_i]_{a_j}|}{|[\omega_i]_{a_j}|}, g^*(\omega_i) = -\log_2(f^*(\omega_i)), \text{ et } h^*(\omega_i) = \sum_{i=1}^n \frac{1}{n} g^*(f^*(\omega_i))$$

Sa version qui tient compte la monotonie, l'entropie de Shannon d'ordre, est définie de la façon suivante :

$$H_S^*(\lambda|a_j) = \sum_{i=1}^n \frac{1}{n} (-\log_2(dsr(\omega_i)))$$

$$\text{où } f^*(\omega_i) = dsr(\omega_i) = \frac{|[\omega_i]_\lambda^{\leq} \cap [\omega_i]_{a_j}^{\leq}|}{|[\omega_i]_{a_j}^{\leq}|}, \text{ et } g^* \text{ et } h^* \text{ sont les mêmes que précédemment.}$$

Les ensembles dominants et classes d'équivalence servant de base pour le calcul de $f^*(\omega_i)$, pour tout $\omega_i \in \Omega$, le stockage de ces structures joue un rôle important dans la complexité du critère de partitionnement.

Étant donnés $\Omega = \{\omega_1, \dots, \omega_n\}$, $\mathcal{A} = \{a_1, \dots, a_m\}$, λ , et $j \in \{1, \dots, m\}$ fixé, on stocke les ensembles dominants générés par l'attribut a_j (respectivement la fonction d'étiquetage λ) dans une matrice notée A (respectivement D) de taille $n \times n$. Pour tout $i, h \in \{1, \dots, n\}$,

- $A_{ih} = 1$ si et seulement si $a_j(\omega_i) \leq a_j(\omega_h)$;
- $D_{ih} = 1$ si et seulement si $\lambda(\omega_i) \leq \lambda(\omega_h)$.

A (respectivement D) est donc rempli en itérant sur chaque $\omega_i \in \Omega$ et en assignant à chaque ligne A_i (respectivement D_i) l'ensemble $\{\omega_h : a_j(\omega_i) \leq a_j(\omega_h)\}$ (respectivement $\{\omega_h : \lambda(\omega_i) \leq \lambda(\omega_h)\}$). Récupérer ce dernier se fait en temps $\mathcal{O}(n)$. On construit donc A et D en temps $\mathcal{O}(n^2)$. Comme les valeurs sont stockées dans une matrice de taille $n \times n$, la complexité spatiale est aussi en $\mathcal{O}(n^2)$.

En ce qui concerne les classes d'équivalence, pour tout $i, h \in \{1, \dots, n\}$,

- $A_{ih} = 1$ si et seulement si $a_j(\omega_i) = a_j(\omega_h)$
- $D_{ih} = 1$ si et seulement si $\lambda(\omega_i) = \lambda(\omega_h)$.

La procédure de construction est la même que celle donnée précédemment, en remplaçant l'inégalité par l'égalité. On obtient les mêmes complexités temporelle et spatiale. Comme les algorithmes manipulant les ensembles dominants et ceux utilisant les classes d'équivalence partagent les mêmes complexités, on étudiera uniquement les fonctions faisant intervenir les ensembles dominants.

Pour a_j et λ fixés, A et D ne sont calculés qu’une seule fois car $[\omega_i]_{a_j}^{\leq}$ et $[\omega_i]_{\lambda}^{\leq}$ restent constants, pour tout $\omega_i \in \Omega$.

On s’intéresse maintenant au calcul de $f^*(\omega_i)$, qu’on appellera également $*-dsr(\omega_i)$. Afin de calculer $dsr(\omega_i)$, on récupère $[\omega_i]_{a_j}^{\leq}$ et $[\omega_i]_{\lambda}^{\leq}$ à l’aide de A et D en temps $\mathcal{O}(n)$. Le calcul de $[\omega_i]_{a_j} \cap [\omega_i]_{\lambda}$ se faisant aussi en $\mathcal{O}(n)$, calculer $dsr(\omega_i)$ se fait en temps $\mathcal{O}(n)$. L’espace occupé est en $\mathcal{O}(n)$. En revanche, le calcul de $maxdsr(\omega_i)$, $mindsr(\omega_i)$ et $avgdsr(\omega_i)$ nécessite d’itérer sur chaque $\omega_h \in [\omega_i]_{a_j}$. Les complexités temporelle et spatiale de ces fonctions sont donc en $\mathcal{O}(n^2)$.

Afin d’observer le comportement des $*-dsr$ sur un attribut monotone, on génère une base de 100 points à deux classes, deux dimensions dont une seule est monotone par rapport à la classe. Ici, la valeur de l’attribut monotone et celle de la classe croissent en fonction de l’indice de l’élément. On duplique et bruitte la même base de sorte à obtenir des bases bruitées à 0%, 25%, 50%, 75%. Les courbes 1, 2, 3, et 4 tracent la valeur de $dsr(\omega_i)$ en fonction de i pour chaque base.

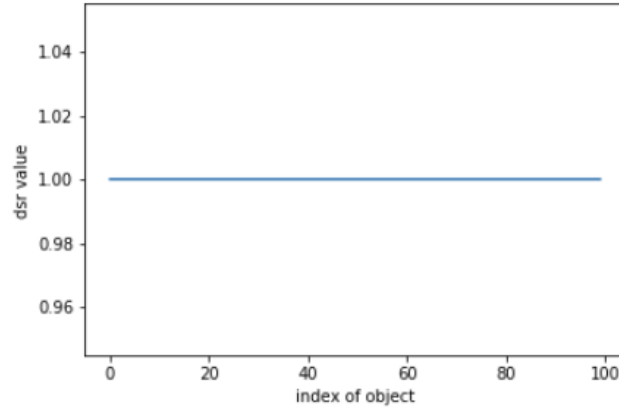


FIGURE 1 – Valeur de dsr en fonction de l’indice de l’objet pour une base non-bruitée

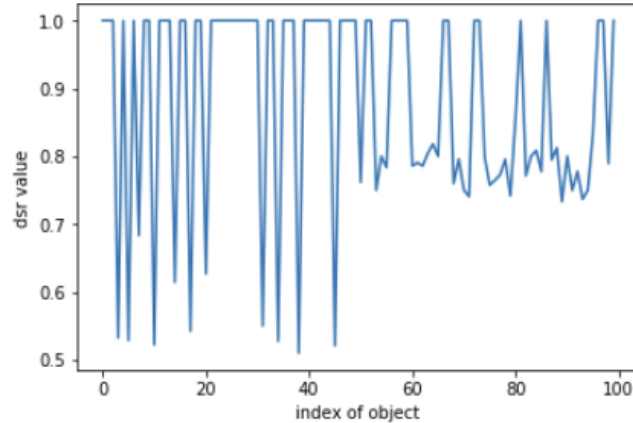


FIGURE 2 – Valeur de dsr en fonction de l’indice de l’objet pour une base bruitée à 25%

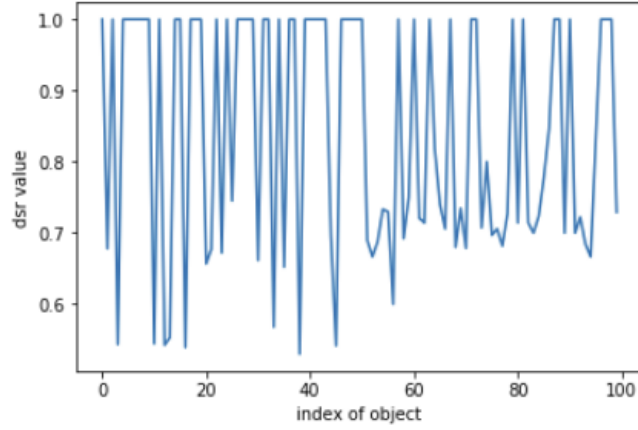


FIGURE 3 – Valeur de dsr en fonction de l'indice de l'objet pour une base bruitée à 50%

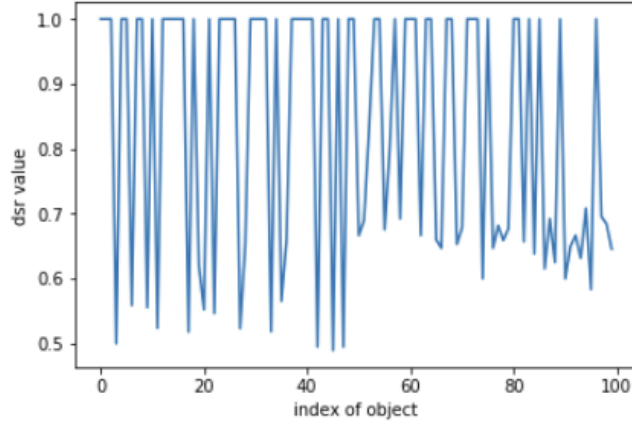


FIGURE 4 – Valeur de dsr en fonction de l'indice de l'objet pour une base bruitée à 75%

Dans la base non bruitée, $\text{dsr}(\omega_i) = 1$ pour chaque $\omega_i \in \Omega$. En effet, on a $a_j(\omega_i) \leq a_j(\omega_h) \Rightarrow \lambda(\omega_i) \leq \lambda(\omega_h)$ si et seulement si $\{\omega_h \in \Omega : \lambda(\omega_i) \leq \lambda(\omega_h) \wedge a_j(\omega_i) \leq a_j(\omega_h)\} = \{\omega_h \in \Omega : a_j(\omega_i) \leq a_j(\omega_h)\}$ si et seulement si $\frac{|\{\omega_i\}_{\lambda}^{\leq} \cap \{\omega_i\}_{a_j}^{\leq}|}{|\{\omega_i\}_{a_j}^{\leq}|} = 1$.

En revanche, plus le taux de bruit dans la base augmente, plus on observe de "pics". En effet, on observe davantage de ω_i tels que $\{\omega_h \in \Omega : a_j(\omega_i) \leq a_j(\omega_h)\}$ est modifié, et donc $|\{\omega_h \in \Omega : \lambda(\omega_i) \leq \lambda(\omega_h) \wedge a_j(\omega_i) \leq a_j(\omega_h)\}| \leq |\{\omega_h \in \Omega : a_j(\omega_i) \leq a_j(\omega_h)\}|$. Les courbes des autres *-dsr sont données en annexe.

Toutes les fonctions g^* présentées dans l'article ont des complexités temporelle et spatiale en $\mathcal{O}(1)$. En revanche, comme les fonctions h^* agrègent les couches g^* correspondant à chaque élément de Ω , elles sont toutes en $\mathcal{O}(n)$.

Le calcul de $H^*(\lambda|a_j)$ se fait donc en :

- $\mathcal{O}(n^2)$ lorsque $f^* = \text{dsr}$
- $\mathcal{O}(n^3)$ lorsque $f^* \in \{\text{mindsr}, \text{maxdsr}, \text{avgdsr}\}$

On étudie, pour chaque couple de mesures de discrimination (H, H') données dans [18], l'évolution de H' en fonction de H pour 2, 3 et 5 classes.

Pour cela, on génère aléatoirement, pour chaque nombre de classes, 100 bases de 100 exemples à un attribut monotone par rapport à la classe. Pour chaque base, on récupère les seuils de coupure engendrés par la discrétisation de l'attribut (étape décrite dans la section suivante) et, pour chaque mesure, on enregistre les valeurs obtenues pour chaque seuil.

Les figures 5, 6, 7 tracent les corrélations entre chaque mesure de discrimination.

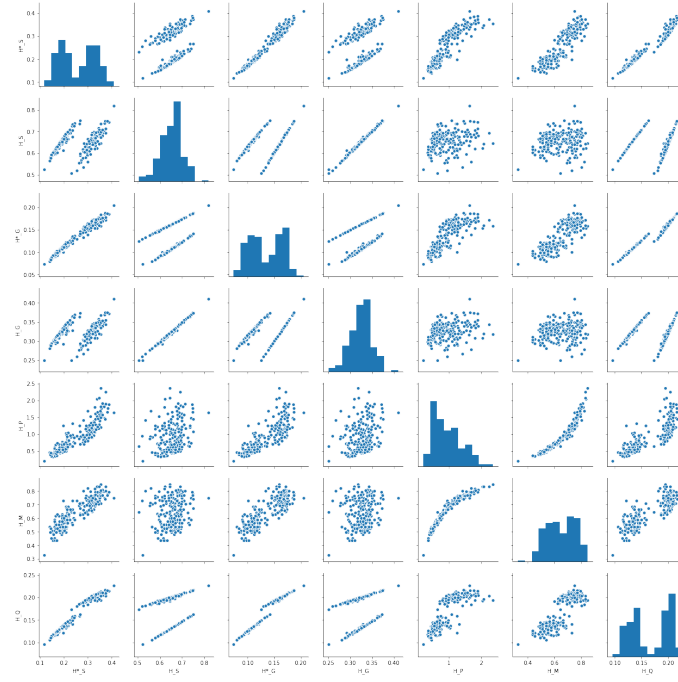


FIGURE 5 – H' en fonction de H (2 classes)

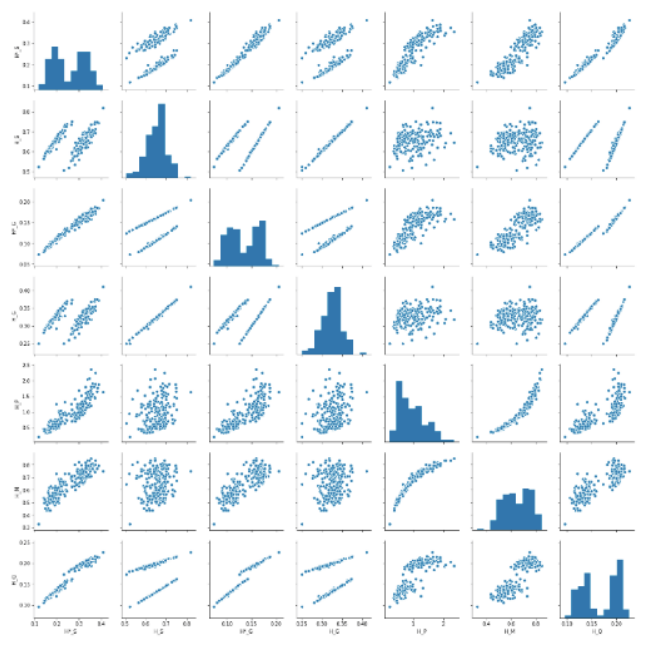


FIGURE 6 – H' en fonction de H (3 classes)

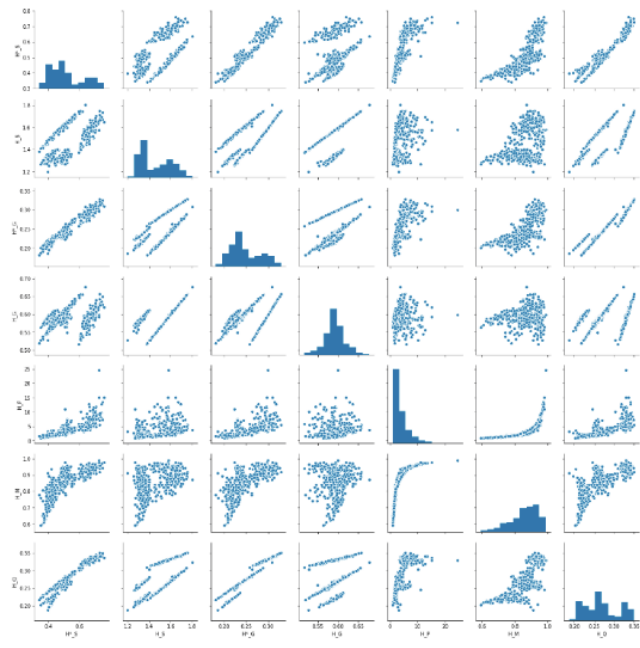


FIGURE 7 – H' en fonction de H (5 classes)

Pour chaque cas, on remarque que les mesures suivantes sont corrélées linéairement :

- H_S^* et H_G^*
- H_G et H_S

Et les mesures suivantes sont corrélées quadratiquement :

- H_P et H_M
- H_Q et H_S^*

4.2 Algorithme de construction d'arbres de décision paramétré par une mesure de discrimination

En reprenant essentiellement RDMT, l'algorithme 2 est un classifieur paramétré par une mesure de discrimination H (d'ordre ou non) et par 3 critères d'arrêt. Comme RDMT, cet algorithme ne produit pas un arbre globalement monotone. Néanmoins, si les données sont monotone-consistantes et si H tient compte de la relation graduelle entre les valeurs d'attributs et les valeurs de classe, alors cette méthode garantit une forme plus faible de monotonie (appelée *rule monotonicity*).

RDMT et notre algorithme diffèrent sur plusieurs points :

- On ne traite pas les attributs non numériques
- Lorsqu'une feuille est créée, on lui assigne la classe majoritaire parmi les exemples ayant servi à sa construction
- On distingue deux mesures de discrimination : l'une doit être minimisée pour le partitionnement et l'autre permet de déterminer l'arrêt. Dans la suite, la deuxième mesure utilisée sera toujours H_S .

Comme l'on se restreint à des arbres de décision binaires, les données doivent être partitionnées en deux sous-ensembles à chaque étape de l'induction. Cette partition se fait en divisant l'ensemble courant Ω_α selon l'attribut a^* qui respecte le plus la contrainte de monotonie, i.e. celui dont la valeur $x_{j_s}^*$ minimise $H(\lambda|a_j^{x_{j_s}^*})$, pour $j = 1, \dots, m, s = 1, \dots, t_j - 1$.

Etant donné une mesure de discrimination H^* , *DISCRETIZE* (algorithme 1) permet de récupérer, pour chaque $a_j \in \mathcal{A}$, le seuil de coupure minimisant la valeur de H ainsi que cette dernière. On parcourt tous les seuils de coupure x_{j_s} de a_j pour trouver x_j^* minimisant $H(\lambda|a_j^{x_{j_s}^*})$, soit :

$$x_j^* = \arg \min \{H(\lambda|a_j^{x_{j_s}^*}), s = 1, \dots, t_j - 1\}$$

Algorithm 1 Discrétisation

procedure DISCRETIZE(H^*, Ω) $\triangleright H^*$: mesure de discrimination construite de façon hiérarchique $\triangleright \Omega$: base de données étiquetées $\triangleright a_j$: attribut à discrétiser $n \leftarrow |\Omega|$ $T \leftarrow$ matrice dont la première colonne contient les valeurs de $a_j(\omega_i)$ triées par ordre croissant, la deuxième, les valeurs de $\lambda(\omega_i)$ triées selon $a_j(\omega_i)$, et la troisième, les $i = 1, \dots, n$ triés selon $a_j(\omega_i)$, pour tout $\omega_i \in \Omega$. $U \leftarrow \{a_j(\omega_i) : \forall \omega_i \in \Omega\}$ \triangleright valeurs uniques de $a_j(\omega_i)$ pour tout $\omega_i \in \Omega$ $a_j^b(\omega_i) \leftarrow 1, \forall \omega_i \in \Omega \quad \triangleright$ valeur de l'attribut a_j binarisée i.e $a_j^b(\omega_i) = 0$ si $a_j(\omega_i) \leq x_{j_s}, 1$ sinon, pour x_{j_s} fixé $[\omega_i]_{a_j^b}^{\leq} \leftarrow \Omega, \forall \omega_i \in \Omega$ $[\omega_i]_{a_i}^{\leq} \leftarrow \{\omega_h \in \Omega : \lambda(\omega_i) \leq \lambda(\omega_h)\}$ $S \leftarrow \emptyset$ \triangleright seuils de coupure considérés $E \leftarrow \emptyset$ \triangleright valeurs de H^* obtenues pour chaque seuil de coupure $s \in S$ $V \leftarrow \emptyset$ \triangleright ensemble des $\omega \in \Omega$ déjà visités**for** $v \in U$ **do** \triangleright on considère chaque valeur unique de \mathcal{A} $\Omega_v \leftarrow \{\omega_i \in \Omega : a_j(\omega_i) = v\}$ $C_v \leftarrow \{c \in C : \exists \omega_h \in \Omega_v, \lambda(\omega_h) = c\}$ $v' \leftarrow$ valeur suivante dans U $\Omega_{v'} \leftarrow \{\omega_i \in \Omega : a_j(\omega_i) = v'\}$ $C_{v'} \leftarrow \{c \in C : \exists \omega_h \in \Omega_{v'}, \lambda(\omega_h) = c\}$ $V \leftarrow V \cup \{\Omega_v\}$ $a_j^b(\omega_i) \leftarrow 0$ pour tout $\omega_i \in \Omega_v$ $x_{j_s} \leftarrow \frac{v+v'}{2}$ **if** $C_{v'} \neq C_v$ **then** $[\omega_i]_{a_j^b} \leftarrow \Omega, \forall \omega_i \in V$ $\bar{V} \leftarrow \Omega \setminus V$ \triangleright ensemble des $\omega \in \Omega$ non visités $[\omega_i]_{a_j^b}^{\leq} \leftarrow \bar{V}, \forall \omega_i \in \bar{V}$ $S \leftarrow S \cup \{x_{j_s}\}$ $E \leftarrow E \cup \{H^*(\lambda|a_j^b)\}$ **return** $\min E, \arg \min_{s \in S} E$ \triangleright retourner la valeur d'entropie minimale et le seuil de coupure permettant de l'obtenir

Algorithm 2 Construction de l'arbre

procedure BUILD_DT($\Omega_\alpha, H^*, H, \epsilon, \Delta, n_{min}, \lambda, \delta$)

▷ Ω_α : base de données étiquetées
▷ H^* : mesure de discrimination utilisée pour le partitionnement
▷ H : mesure de discrimination utilisée pour déterminer l'arrêt
▷ ϵ : limite inférieure pour H
▷ Δ : longueur maximale d'un chemin de la racine à une feuille
▷ n_{min} : taille minimale pour Ω_α
▷ λ : fonction d'étiquetage
▷ δ : longueur du chemin de la racine au noeud courant

$h \leftarrow H(\Omega_\alpha, \lambda)$

if $h < \epsilon$ **or** $|\Omega_\alpha| < n_{min}$ **or** $\forall \omega_i, \omega_h \in \Omega_\alpha, \lambda(\omega_i) = \lambda(\omega_h)$ **or** $\delta > \Delta$ **then**

$c_\alpha \leftarrow c$ tel que $|\{\omega_i \in \Omega_\alpha : \lambda(\omega_i) = c\}| \geq |\{\omega_i \in \Omega_\alpha : \lambda(\omega_i) \neq c\}|$

return LEAF(c_α, Ω)

$m \leftarrow$ nombre d'attributs dans Ω

$S \leftarrow \emptyset$ ▷ pour chaque $a_j \in \mathcal{A}$, contient $x_{j*} = \arg \min\{H(\lambda|a_j^{x_{js}}), s = 1, \dots, t_j - 1\}$

$E \leftarrow \emptyset$ ▷ pour chaque $a_j \in \mathcal{A}$, contient $\min\{H(\lambda|a_j^{x_{js}}), s = 1, \dots, t_j - 1\}$

for $a_j=0$ to $m-1$ **do**

if $a_j(\omega_i) = a_j(\omega_h), \forall \omega_i, \omega_h \in \Omega_\alpha$ **then**

$S \leftarrow S \cup \{\infty\}$

$E \leftarrow E \cup \{\infty\}$

else

$x_{j*}, h \leftarrow \text{DISCRETIZE}(H^*, \Omega_\alpha, a_j)$

$S \leftarrow S \cup \{x_{j*}\}$

$E \leftarrow E \cup \{h\}$

$x_* \leftarrow \arg \min_{s \in S} E$

$a_*^{x_*} \leftarrow \arg \min_{a_j \in \mathcal{A}} E$

$\Omega_{\leq}, \Omega_{\geq} \leftarrow \text{DIVIDE}(\Omega_\alpha, a_*^{x_*}, x_*)$

if $|\Omega_{\leq}| = 0$ **then**

return $\{c\}$ tel que $|\{\omega_i \in \Omega_{\geq} : \lambda(\omega_i) = c\}| \geq |\{\omega_i \in \Omega_{\geq} : \lambda(\omega_i) \neq c\}|$

if $|\Omega_{\geq}| = 0$ **then**

return $\{c\}$ tel que $|\{\omega_i \in \Omega_{\leq} : \lambda(\omega_i) = c\}| \geq |\{\omega_i \in \Omega_{\leq} : \lambda(\omega_i) \neq c\}|$

$\mathcal{T}_{\leq} \leftarrow \text{BUILD_DT}(\Omega_{\leq}, H^*, H, \epsilon, \Delta, n_{min}, \lambda, \delta + 1)$

$\mathcal{T}_{\geq} \leftarrow \text{BUILD_DT}(\Omega_{\geq}, H^*, H, \epsilon, \Delta, n_{min}, \lambda, \delta + 1)$

$\mathcal{T} \leftarrow \text{TREE}(\{a_*^{x_*}\}, \mathcal{T}_{\leq}, \mathcal{T}_{\geq})$

return \mathcal{T}

Considérons la base de la figure 8 contenant 60 points à deux dimensions x et y et à 3 classes. x est l'attribut dont les valeurs sont monotones par rapport aux valeurs de la classe. Chaque classe est représentée par le même nombre de points. La classe 1 est représentée par les points bleus, la 2 par les points oranges, et la 3 par les points verts.

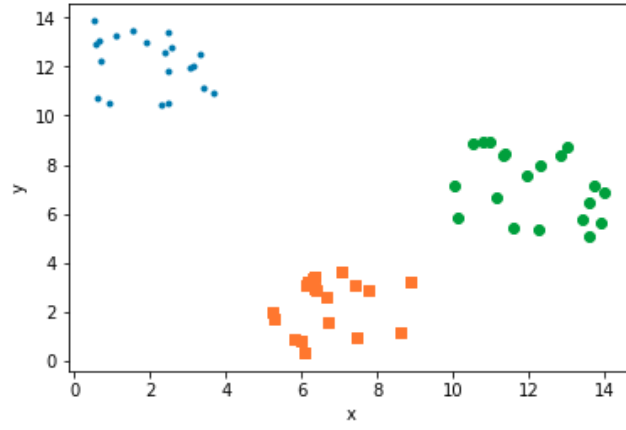
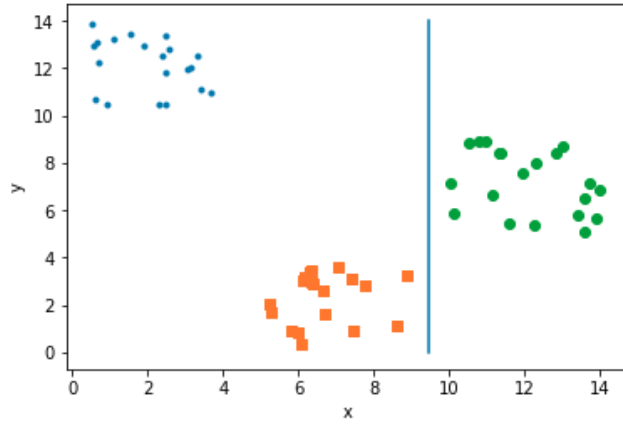
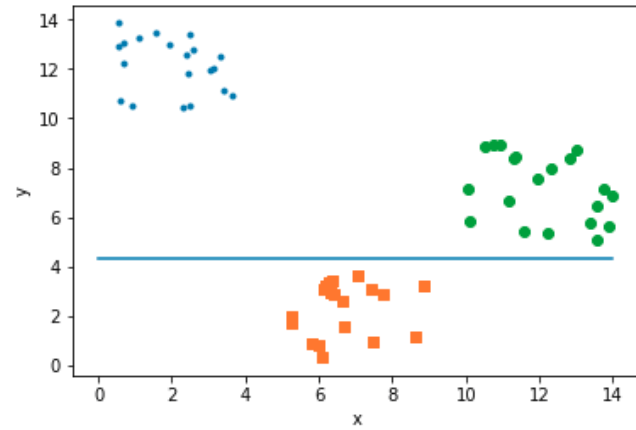


FIGURE 8 – Base de données jouet



$$H_S^*(\lambda|\mathbf{x}) = 0.19 \leq H_S^*(\lambda|\mathbf{y}) = 0.53$$

FIGURE 9 – Seuil de coupure généré par H_S^*



$$H_S(\lambda|\mathbf{x}) = H_S(\lambda|\mathbf{y}) = 0.67$$

FIGURE 10 – Seuil de coupure généré par H_S

Dataset	Nombre d'exemples	Nombre d'attributs	Nombre de classes
Breast Cancer	286	8	2
Cars Evaluation	1728	6	4
CPU	209	9	8

TABLE 1 – Description des bases de classification

Dataset		H_S^*	H_S	H_G^*	H_G	H_P^*	H_M^*	H_O^*
CPU	accuracy	68.40 % \pm 0.04 %	66.02 % \pm 0.02 %	66.48 % \pm 0.05 %	66.97 % \pm 0.04 %	66.51 % \pm 0.00 %	60.73 % \pm 0.07 %	67.93 % \pm 0.02 %
	profondeur	10.50 \pm 0.50	9.00 \pm 0.00	10.00 \pm 1.00	8.50 \pm 0.50	11.00 \pm 0.00	14.50 \pm 0.50	11.00 \pm 1.00
	feuilles	35.00 \pm 2.00	27.50 \pm 3.50	36.50 \pm 2.50	28.50 \pm 4.50	36.50 \pm 0.50	41.50 \pm 1.50	40.50 \pm 3.50
	ratio	77.99 % \pm 0.11 %	86.00 % \pm 2.83 %	80.86 % \pm 1.05 %	81.80 % \pm 0.41 %	73.23 % \pm 4.66 %	70.73 % \pm 2.35 %	50.64 % \pm 16.03 %
	paires de feuilles	9.50 \pm 0.50	7.50 \pm 1.50	9.50 \pm 0.50	7.50 \pm 1.50	9.50 \pm 0.50	10.50 \pm 0.50	12.50 \pm 1.50
Breast Cancer	accuracy	59.17 % \pm 0.13 %	63.77 % \pm 0.11 %	56.31 % \pm 0.14 %	62.03 % \pm 0.08 %	61.40 % \pm 0.13 %	59.80 % \pm 0.11 %	54.80 % \pm 0.13 %
	profondeur	20.00 \pm 1.48	16.60 \pm 1.74	20.40 \pm 1.36	15.10 \pm 1.45	19.70 \pm 2.24	20.80 \pm 1.60	20.10 \pm 1.58
	feuilles	129.20 \pm 4.71	81.60 \pm 5.99	129.80 \pm 4.64	80.10 \pm 4.25	114.90 \pm 5.34	127.10 \pm 5.45	134.90 \pm 5.66
	ratio	84.51 % \pm 2.87 %	80.15 % \pm 4.76 %	84.49 % \pm 3.07 %	78.44 % \pm 3.88 %	86.80 % \pm 2.86 %	83.67 % \pm 1.80 %	81.68 % \pm 2.81 %
	paires de feuilles	31.90 \pm 2.91	18.90 \pm 1.81	31.60 \pm 3.10	20.60 \pm 1.69	28.30 \pm 2.00	32.40 \pm 1.69	36.30 \pm 2.61
Cars Evaluation	accuracy	81.94 % \pm 0.04 %	81.85 % \pm 0.06 %	82.00 % \pm 0.04 %	78.72 % \pm 0.09 %	87.08 % \pm 0.04 %	88.59 % \pm 0.07 %	78.53 % \pm 0.03 %
	profondeur	10.50 \pm 0.50	10.75 \pm 0.43	10.50 \pm 0.50	10.25 \pm 0.43	10.00 \pm 0.00	11.00 \pm 0.00	11.00 \pm 0.00
	feuilles	66.50 \pm 13.16	37.25 \pm 1.30	66.50 \pm 13.16	38.00 \pm 1.58	52.00 \pm 3.94	54.00 \pm 3.74	121.25 \pm 14.04
	ratio	82.17 % \pm 5.00 %	83.82 % \pm 1.39 %	82.76 % \pm 5.56 %	84.68 % \pm 1.84 %	79.25 % \pm 2.93 %	77.33 % \pm 2.90 %	52.19 % \pm 4.34 %
	paires de feuilles	27.75 \pm 7.36	11.75 \pm 1.09	27.75 \pm 7.36	12.00 \pm 1.41	18.25 \pm 1.92	18.00 \pm 2.12	46.75 \pm 8.79

TABLE 2 – Résultats sur les tâches de classification

On observe sur les figures 9 et 10 les seuils de coupures générés respectivement par H_S^* et H_S sur cette base jouet. On remarque que H_S^* tient compte de la monotonie des valeurs de la classe par rapport aux valeurs de l'attribut x : la valeur d'entropie minimale calculée pour l'attribut x est inférieure à celle calculée pour l'attribut y. Ce n'est cependant pas le cas pour H_S qui est insensible à la monotonie : la coupure optimale sur l'axe x génère la même valeur d'entropie que celle faite sur l'axe y.

4.3 Analyse expérimentale

Dans cette section, on cherche à comparer entre eux les arbres obtenus en utilisant H_S^* , H_S , H_G^* , et H_G comme mesure de partitionnement. L'objectif est d'évaluer l'impact du critère de partitionnement sur les performances des classifieurs en termes de taux de bonne classification et de monotonie. On étudie également l'évolution de ces performances en fonction du taux de bruit non-monotone.

Dans les expériences menées, les arbres sont construits avec l'algorithme 2, en faisant varier la mesure de discrimination parmi H_S^* , H_S , H_G^* , et H_G . Ils sont testés sur des bases artificielles et réelles.

4.3.1 Expérimentation sur des bases réelles

Dans cette partie, nos arbres sont testés sur différentes bases réelles de classification et une base de régression récupérées sur UCI [8]. Pour chaque base, une étape de pré-traitement est effectuée : les attributs ordonnés mais catégoriels sont numérisés, et les attributs ne présentant pas d'ordre ainsi que les exemples incomplets sont supprimés.

Les arbres sont évalués sur leur accuracy, leur profondeur, leur nombre de feuilles, leur ratio de paires de feuilles non-monotones, ainsi que le nombre de paires de feuilles.

On cherche à déterminer si, pour chaque mesure d'évaluation, les mesures de discrimination testées produisent des résultats significativement différents.

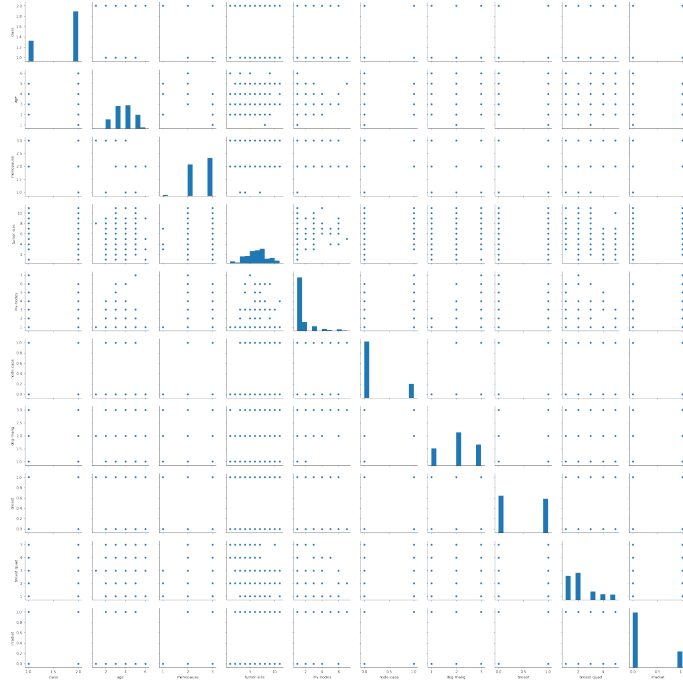


FIGURE 11 – Corrélations entre les attributs de la base *Breast Cancer* et la classe

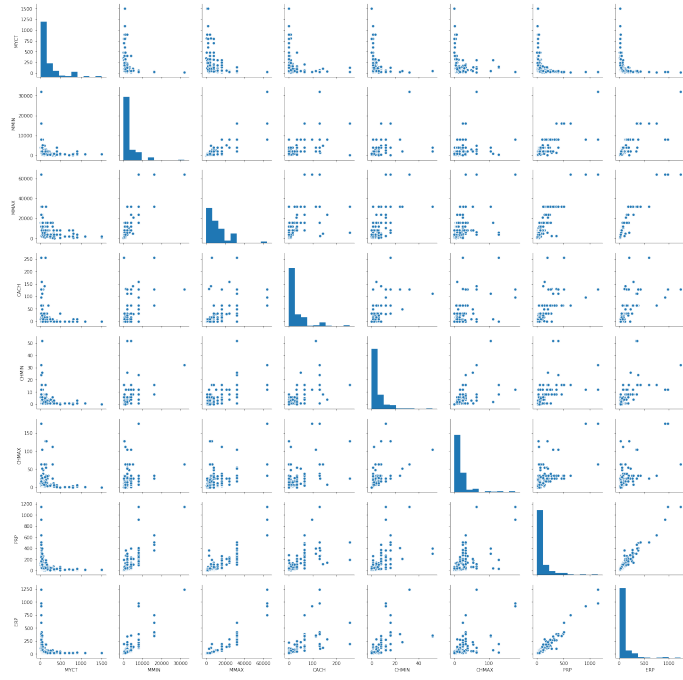


FIGURE 12 – Corrélations entre les attributs de la base *CPU* et la classe

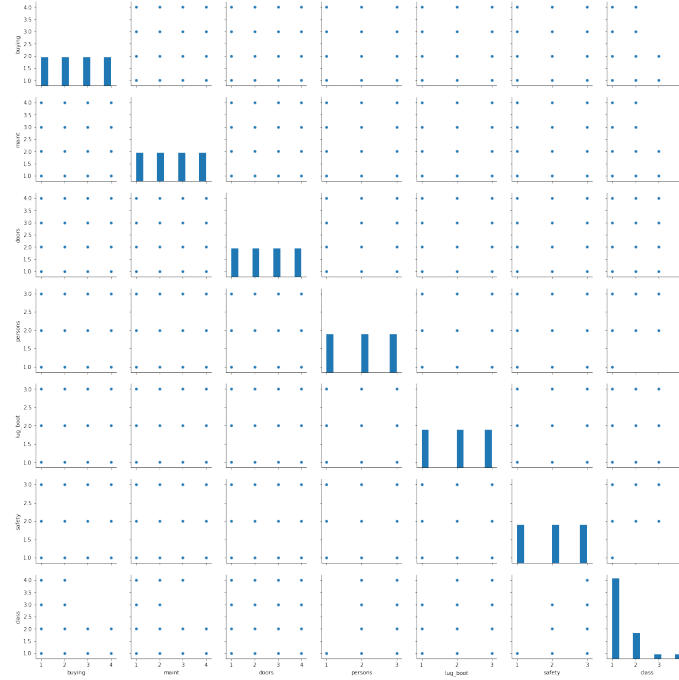


FIGURE 13 – Corrélations entre les attributs de la base *Cars Evaluation* et la classe

Les bases de classification sont décrites dans le tableau 1 p. 17. Le tableau 2 p. 17 contient les résultats des tests menés sur ces bases.

De manière générale, on peut observer que les arbres construits à partir des mesures de discrimination d'ordre sont plus profonds, contiennent plus de feuilles et de paires de feuilles que ceux générés par les mesures d'entropie classiques.

La base *Breast Cancer* comporte 9 attributs dont le nombre de valeurs différentes varie de 2 à 11. Sur cette base, les mesures d'entropie classique produisent de meilleurs résultats en termes de taux de bonne classification et de monotonie que les mesures de discrimination d'ordre. La prise en compte de la monotonie des données dans le critère de partitionnement n'offre donc pas d'avantage. Néanmoins, cette base ne contient que 2 valeurs pour la classe, et les attributs comportent peu de valeurs différentes : entre 2 et 11. Il n'y a donc pas de véritable relation graduelle entre les attributs et la classe. La figure 11 p. 18 confirme l'absence de corrélation entre les valeurs d'attributs et de classe : il n'existe pas de contrainte de monotonie dans cette base.

En revanche, H_S^* , H_G^* , et H_Q^* produisent de meilleurs taux de bonne classification que les mesures d'entropie classiques sur la base *CPU*. De plus, sur les données de *Cars Evaluation*, H_S^* , H_G^* , H_P^* et H_M^* génèrent de meilleures prédictions que les mesures classiques. Sur ces deux bases, tous les ratio produits par les mesures de discrimination d'ordre sont plus faibles que ceux produits par les mesures classiques : cela indique que les arbres générés par une mesure qui prend en compte la relation graduelle entre les attributs et la classe garantissent une plus forte monotonie que les autres.

Dans le but de comprendre pourquoi les mesures de discrimination d'ordre présentent de meilleures performances, on examine les corrélations entre chaque paire d'attributs, classe incluse (*PRP*, avant-dernière colonne en partant de la gauche, avant-dernière ligne en partant du haut). La figure 12 p. 18 montre que, dans la base *CPU*, les valeurs des attributs sont très étalées et un certain nombre d'attributs apparaissent corrélés (positivement comme *ERP* et *PRP*, ou négativement comme *PRP* et *MYCT*).

Cependant, on peut voir sur la figure 13 p. 19 que les attributs de la base *Cars Evaluation* comprennent peu de valeurs différentes et sont peu corrélés à la classe (*class*, dernière colonne en partant de la gauche, dernière ligne en partant du haut) : on manque encore d'une explication générale.

Nombre de classes		H_S^*	H_S	H_G^*	H_G
2 classes	accuracy	68.93% \pm 0.11 %	69.68 % \pm 0.11 %	69.54 % \pm 0.12 %	70.33% \pm 0.12 %
	profondeur	17.0 \pm 1.22	13.5 \pm 0.87	17.5 \pm 1.50	12.25 \pm 0.43
	feuilles	99.75 \pm 5.07	78.25 \pm 5.54	102.5 \pm 7.79	79.25 \pm 5.67
	ratio	87.82% \pm 0.01 %	87.86% \pm 0.01 %	88.36% \pm 0.01 %	87.42% \pm 0.02 %
	paires de feuilles	28.5 \pm 2.29	24.5 \pm 2.87	28.5 \pm 2.60	29.0 \pm 1.22
3 classes	accuracy	39.25% \pm 0.10 %	33.00% \pm 0.12%	35.14 % \pm 0.07%	32.33 % \pm 0.11%
	profondeur	20.5 \pm 1.66	14.75 \pm 1.09	18.25 \pm 1.48	15.5 \pm 2.69
	feuilles	158.75 \pm 5.12	140.5 \pm 1.66	164.25 \pm 3.83	143.75 \pm 5.72
	ratio	74.82 % \pm 0.03%	80.74 % \pm 0.04%	75.60 % \pm 0.04%	79.11 % \pm 0.04%
	paires de feuilles	44.5 \pm 2.60	48.25 \pm 1.92	45.75 \pm 3.03	47.25 \pm 3.63
7 classes	accuracy	31.18 % \pm 0.08 %	30.37 % \pm 0.06%	29.37 % \pm 0.08 %%	29.24% % \pm 0.05 %
	profondeur	19.0 \pm 1.87	15.75 \pm 1.29	19.5 \pm 1.12	14.25 \pm 1.30
	feuilles	171.75 \pm 3.11	151.5 \pm 3.77	172.5 \pm 3.64	157.5 \pm 4.09
	ratio	53.82 % \pm 0.05 %	56.23 % \pm 0.02 %	51.74% % \pm 0.01 %	51.89 % \pm 0.05 %
	paires de feuilles	47.0 \pm 2.23	48.25 \pm 1.92	47.0 \pm 1.87	51.25 \pm 3.34
9 classes	accuracy	25.60 % \pm 0.05 %	23.93 % \pm 0.08%	24.53 % \pm 0.06 %	22.07% % \pm 0.05 %
	profondeur	19.5 \pm 1.5	14.75 \pm 0.43	20.0 \pm 1.22	14.0 \pm 1.00
	feuilles	171.75 \pm 5.89	158.5 \pm 0.87	176.0 \pm 3.94	160.0 \pm 2.55
	ratio	53.07 % \pm 0.04 %	55.45 % \pm 0.03 %	51.86% % \pm 0.02 %	53.41% % \pm 0.01 %
	paires de feuilles	46.0 \pm 4.53	52.25 \pm 3.69	48.0 \pm 1.87	54.0 \pm 2.12

TABLE 3 – Résultats sur la base de régression



FIGURE 14 – Corrélations entre les attributs de la base *Airfoil Self-Noise* et la classe

D'autre part, la base de régression *Airfoil Self-Noise* est utilisée afin d'étudier l'impact du nombre de classes sur les performances du classifieur. Cette base est constituée de 1503 éléments à 5 attributs : *frequency* (21 valeurs différentes), *angle of attack* (27 valeurs différentes), *chord length* (6 valeurs différentes), *velocity* (4 valeurs différentes), et *thickness* (105 valeurs différentes). La classe à prédire, *scaled sound pressure level*, comporte 1456 valeurs différentes. En fonction des quantiles, on discrétise les valeurs de classe en 2, 3, 7, et 9 classes afin d'obtenir des tâches de classification. On étudie les performances du classifieur en fonction de la mesure de discrimination utilisée pour chaque discrétisation de classe, puis l'impact du nombre de classes sur les résultats.

Comme précédemment, les arbres générés par des mesures de discrimination d'ordre sont plus profonds et contiennent plus de feuilles que les autres, pour toute discrétisation. De plus, ils permettent d'obtenir de meilleurs résultats en termes de taux de bonne classification et de monotonie. En revanche, lorsque la base contient plus de 2 classes, les arbres construits par des mesures d'en-

trope classiques comprennent plus de paires de feuilles que ceux tenant compte de la monotonie.

On remarque également, pour chaque mesure étudiée, une nette diminution du ratio de paires non-monotones lorsque l'on augmente le nombre de classes. Plus il est élevé, plus la relation graduelle existant entre les attributs et la classe est renforcée. Couplée par des mesures de discrimination prenant en compte la monotonie, on obtient des classifieurs plus performants que ceux générés par des mesures d'entropie classiques.

Néanmoins, la figure 14 p. 20 montre le peu de corrélation entre les attributs et la classe (représentée par la dernière ligne en partant du haut, et la dernière colonne en partant de la gauche). Il est encore difficile d'expliquer cette différence de performance entre les classifieurs monotones et non-monotones.

Les expériences menées sur les bases précédentes mettent en valeur les avantages des arbres générés par des mesures tenant compte de la contrainte de monotonie par rapport aux arbres construits avec des mesures classiques. Cependant, on manque d'une explication générale à ce phénomène : l'absence de relation monotone entre les attributs et la classe n'implique pas de moins bonnes performances de la part des arbres monotones par rapport aux arbres classiques. De plus, comme pour les expériences menées sur des données artificielles, il est toujours difficile d'évaluer le taux de monotonie d'un arbre.

5 Conclusion du stage et perspectives

Mon stage a porté sur un framework particulier de classification monotone [18]. Après avoir fait l'état de l'art sur les approches de classification monotone proposées dans la littérature, je me suis placée dans le contexte des arbres de décision monotones paramétrés par une mesure de sélection d'attribut.

J'ai effectué une étude théorique des mesures de discrimination d'ordre afin de vérifier certaines propriétés clés et d'en extraire d'autres (notamment le sens de variation).

J'ai également implémenté une bibliothèque pour les arbres monotones comprenant le modèle de construction hiérarchique des mesures [18], ainsi que l'algorithme de construction d'arbres monotones [18] légèrement modifié.

Des expérimentations sur des bases artificielles et réelles ont été menées afin de mettre en avant les avantages des arbres monotones et comparer les différentes mesures de discrimination entre elles. Les résultats observés sur les bases réelles dans lesquelles il existe une forte gradualité entre les attributs et la classe montrent que les arbres produits par des mesures de discrimination d'ordre sont plus performants et plus monotones que ceux générés par des mesures d'entropie classiques. Cependant, les tests lancés sur des bases générées artificiellement ne permettent pas de souligner l'apport des mesures de discrimination d'ordre en termes de monotonie des arbres.

Une perspective possible serait de travailler sur ce problème : il peut venir de la façon de partitionner les données (discrétisation en deux sous-ensembles alors qu'il peut y avoir plus de deux classes) ou de la mesure d'évaluation de la monotonie utilisée. De plus, notre algorithme considère que tous les attributs sont monotones par rapport à la classe, ce qui n'est cependant pas le cas dans la plupart des cas d'applications réelles.

Sur un plan plus personnel, ce stage m'a initiée et donné goût à la recherche en laboratoire. Je tiens à remercier mon maître de stage, Christophe Marsala, ainsi qu'Arthur Guillon, doctorant à LFI, pour m'avoir aidée et accompagnée durant ce stage.

Références

- [1] A. BEN-DAVID. "Monotonicity maintenance in information-theoretic machine learning algorithms". In : *Machine Learning* 19.1 (1995), p. 29–43.
- [2] A. BEN-DAVID, L. STERLING et Y.-H. PAO. "Learning and classification of monotonic ordinal concepts". In : *Computational Intelligence* 5.1 (1989), p. 45–49.
- [3] A. BEN-DAVID, L. STERLING et T. TRAN. "Adding monotonicity to learning algorithms may impair their accuracy". In : *Expert Systems with Applications* 36.3 (2009), p. 6627–6634.

- [4] J. BŁASZCZYŃSKI, R. SŁOWIŃSKI et M. SZELĄG. “Induction of ordinal classification rules from incomplete data”. In : *International Conference on Rough Sets and Current Trends in Computing*. Springer. 2012, p. 56–65.
- [5] J. BŁASZCZYŃSKI, R. SŁOWIŃSKI et M. SZELĄG. “Sequential covering rule induction algorithm for variable consistency rough set approaches”. In : *Information Sciences* 181.5 (2011), p. 987–1002.
- [6] K. CAO-VAN. “Supervised ranking : from semantics to algorithms”. Thèse de doct. Ghent University, 2003.
- [7] K. CAO-VAN et B. DE BAETS. “Consistent representation of rankings”. In : *Theory and Applications of Relational Structures as Knowledge Instruments*. Springer, 2003, p. 107–123.
- [8] D. DHEERU et E. KARRA TANISKIDOU. *UCI Machine Learning Repository*. 2017.
- [9] S. GRECO, B. MATARAZZO et R. SLOWINSKI. “A new rough set approach to evaluation of bankruptcy risk”. In : *Operational tools in the management of financial risks*. Springer, 1998, p. 121–136.
- [10] S. GRECO, B. MATARAZZO et R. SLOWINSKI. “Customer satisfaction analysis based on rough set approach”. In : *Zeitschrift für Betriebswirtschaft* 77.3 (2007), p. 325–339.
- [11] S. GRECO, B. MATARAZZO et R. SLOWINSKI. “Rough approximation by dominance relations”. In : *International journal of intelligent systems* 17.2 (2002), p. 153–171.
- [12] S. GRECO, B. MATARAZZO et R. SLOWINSKI. “Rough sets methodology for sorting problems in presence of multiple attributes and criteria”. In : *European journal of operational research* 138.2 (2002), p. 247–259.
- [13] Q. HU, X. CHE, L. ZHANG, D. ZHANG, M. GUO et D. YU. “Rank entropy-based decision trees for monotonic classification”. In : *IEEE Transactions on Knowledge and Data Engineering* 24.11 (2012), p. 2052–2064.
- [14] J. D. HUNTER. “Matplotlib : A 2D graphics environment”. In : *Computing In Science & Engineering* 9.3 (2007), p. 90–95.
- [15] W. KOTŁOWSKI et R. SŁOWIŃSKI. “Rule learning with monotonicity constraints”. In : *Proceedings of the 26th Annual International Conference on Machine Learning*. ACM. 2009, p. 537–544.
- [16] B. LEO, J. H. FRIEDMAN, R. A. OLSHEN et C. J. STONE. “Classification and regression trees”. In : *Wadsworth International Group* (1984).
- [17] C. MARSALA. “Gradual fuzzy decision trees to help medical diagnosis”. In : *Fuzzy Systems (FUZZ-IEEE), 2012 IEEE International Conference on*. IEEE. 2012, p. 1–6.
- [18] C. MARSALA et D. PETTURITI. “Rank discrimination measures for enforcing monotonicity in decision tree induction”. In : *Information Sciences* 291 (2015), p. 143–171.
- [19] M. J. PAZZANI, S. MANI et W. R. SHANKLE. “Acceptance of rules generated by machine learning among medical experts”. In : *Methods of information in medicine* 40.05 (2001), p. 380–385.
- [20] F. PEDREGOSA, G. VAROQUAUX, A. GRAMFORT, V. MICHEL, B. THIRION, O. GRISEL, M. BLONDEL, P. PRETTENHOFER, R. WEISS, V. DUBOURG, J. VANDERPLAS, A. PASSOS, D. COURNAPEAU, M. BRUCHER, M. PERROT et E. DUCHESNAY. “Scikit-learn : Machine Learning in Python”. In : *Journal of Machine Learning Research* 12 (2011), p. 2825–2830.
- [21] S. PEI et Q. HU. “Partially monotonic decision trees”. In : *Information Sciences* 424 (2018), p. 104–117.
- [22] Y. QIAN, H. XU, J. LIANG, B. LIU et J. WANG. “Fusing monotonic decision trees”. In : *IEEE Transactions on Knowledge and Data Engineering* 27.10 (2015), p. 2717–2728.
- [23] J. R. QUINLAN. “Induction of decision trees”. In : *Machine learning* 1.1 (1986), p. 81–106.
- [24] S. v. d. WALT, S. C. COLBERT et G. VAROQUAUX. “The NumPy Array : A Structure for Efficient Numerical Computation”. In : *Computing in Science and Engg.* 13.2 (mar. 2011), p. 22–30.

- [25] H. WANG, M. ZHOU et K. SHE. “Induction of ordinal classification rules from decision tables with unknown monotonicity”. In : *European Journal of Operational Research* 242.1 (2015), p. 172–181.
- [26] F. XIA, W. ZHANG, F. LI et Y. YANG. “Ranking with decision tree”. In : *Knowledge and information systems* 17.3 (2008), p. 381–395.