

Méthode basée sur les ensembles approximatifs pour l'apprentissage incrémental en présence des données déséquilibrées

Sarra Bouzayane* ***, Inès Saad* **

*Université de Picardie Jules verne, Amiens
{sarra.bouzayane, ines.saad}@u-picardie.fr

**Ecole supérieure de commerce, Amiens

***Institut Supérieur d'Informatique et de Multimédia, Sfax

Résumé. Ce papier propose une méthode basée sur la théorie des ensembles approximatifs et dédiée à l'apprentissage supervisé incrémental dans un contexte de données déséquilibrées. Cette méthode consiste en trois phases : la construction d'une table de décision, l'inférence d'un ensemble de règles de décision et la classification de chaque action potentielle dans l'une des classes de décision prédéfinies. La méthode *MAI2P* est validée dans le contexte des MOOCs (*Massive Open Online Courses*).

1 Introduction

Lorsque les exemples d'apprentissage sont fournis de manière séquentielle, l'apprentissage incrémental pour une prise de décision s'avère une obligation (Greco et al., 2004). Généralement, la phase d'apprentissage est traitée par les techniques conventionnelles de l'apprentissage machine. Cependant, ces techniques demeurent sensibles au problème des données déséquilibrées qui résulte de la répartition inégale entre les instances des classes de décision. Cette inégalité affecte considérablement la qualité de la décision en particulier quand il s'agit de données massives. Ce problème peut, toutefois, être surmonté par l'approche DRSA (*Dominance-based Rough Set Approach*) (Greco et al., 2001) qui repose sur les préférences et l'expertise des décideurs humains pour la construction d'un ensemble d'apprentissage afin de garantir la répartition égale des instances sur l'ensemble de classes de décision.

Ce travail propose une méthode *MAI2P* (*Multicriteria Approach for the Incremental Periodic Prediction*) basée sur l'approche DRSA pour la classification multicritère incrémentale et périodique. La méthode *MAI2P* se compose de trois phases. La première vise à construire une table de décision et repose sur trois étapes : l'identification d'un ensemble d'apprentissage ; la construction d'une famille cohérente de critères pour la caractérisation des actions ; et la classification de chaque action d'apprentissage dans l'une des classes de décision. La deuxième phase est basée sur notre algorithme *DRSA-Incremental* (Bouzayane et Saad, 2017) pour l'inférence et la mise à jour de l'ensemble de règles de décision. La troisième consiste à la classification des "Actions potentielles", en utilisant les règles précédemment inférées. L'approche *MAI2P* est validée sur le contexte des MOOCs (*Massive Open Online Courses*).

Le papier est structuré comme suit : La section 2 définit les notions de base de l'approche DRSA. La section 3 présente un état de l'art. La section 4 détaille la méthode *MAI2P*. La section 5 discute les résultats de l'expérimentation. La section 6 conclut le papier.

2 Préliminaires : Dominance-based Rough Set Approach

L'approche DRSA développée par Greco et al. (2001) est dédiée au problème de tri en aide multicritère à la décision et inspirée de la théorie des ensembles approximatifs. Elle permet de comparer des actions à travers une relation de dominance, rendant compte des préférences d'un décideur, afin d'inférer les règles de décision. Cette approche définit une table d'information par un 4-uplets $S = \langle A, F, V, f \rangle$ tels que : A est un ensemble fini des actions de référence ; F est une famille cohérente de critères ; V est un ensemble des valeurs possibles des critères ; et $f : A \times F \rightarrow V$ est une fonction d'information tel que $f(x, g) \in V_g, \forall x \in A, \forall g \in F$. Chaque action de référence est affectée à une seule classe $Cl_t; t \in \{1, \dots, N\}$.

Relation de dominance : La relation de dominance D_P est définie comme suit : $\forall (x, y) \in A^2, x D_P y \Leftrightarrow f(x, g_j) \succcurlyeq f(y, g_j) \forall g_j \in P \subseteq F, \forall x \in A$, est associé un ensemble, *P-dominant*, d'actions dominant x et un ensemble, *P-dominé*, d'actions dominées par x .

Union inférieure (supérieure) : $Cl_n^{\leq} = \cup_{s \leq n} Cl_s$ ($Cl_n^{\geq} = \cup_{s \geq n} Cl_s$); $n = \{1 \dots N\}$: L'union inférieure (supérieure) de Cl_n signifie que " x appartient au maximum (minimum) à la classe Cl_n ou bien à une classe au mieux (moins) aussi bonne que Cl_n ".

Approximation inférieure $\underline{P}(Cl_n^{\geq})$ (ou $\underline{P}(Cl_n^{\leq})$) : regroupe toutes les actions dont l'ensemble P-dominant (P-dominé) est affecté avec *certitude* à des classes au moins (mieux) aussi bonnes que Cl_n . En revanche, l'approximation supérieure regroupe toutes les actions dont l'affectation est réalisée d'une *manière possible*.

Règles de décision : L'ensemble de règles de décision est appelé modèle de préférences. Ces règles sont générées à partir de l'approximation inférieure et se présentent sous la forme :

Si $f(x, g_1) \geq r_1 \wedge \dots \wedge f(x, g_n) \geq r_n$ alors $x \in Cl_t^{\geq}$ tel que $(r_1, \dots, r_n) \in (V_{g_1} \times \dots \times V_{g_n})$.

3 Travaux antérieurs

Quelques approches dynamiques ont été proposées dans la littérature pour la mise à jour incrémentale des règles de décision suite à la variation de l'ensemble d'apprentissage.

Les auteurs dans (Greco et al., 2004) ont proposé un algorithme appelé *Glance*. Cet algorithme est basé sur des actions négatives. En effet, chaque règle d'une union donnée doit impérativement ne pas satisfaire x , si x n'appartient pas à cette union, mais elle peut aussi ne satisfaire aucune action x et dans ce cas elle demeure sans supports. Ces règles sans supports sont dites non-robustes et donc l'algorithme est aussi dit non robuste. L'algorithme *Glance* stocke dans la mémoire uniquement les règles de décision et pas les exemples d'apprentissage et donc il est économe par rapport à l'utilisation de l'espace mémoire. La complexité de l'algorithme est linéaire si on considère le nombre d'actions et elle est exponentielle en considérant le nombre de critères. Les auteurs dans (Li et al., 2013) ont proposé un algorithme de mise à jour incrémentale des approximations inférieures et supérieures de l'approche DRSA lors de l'ajout (ou la suppression) d'une seule action dans le système d'information. La méthode

nécessite : premièrement, la mise à jour des unions inférieures et supérieures des classes de décision, deuxièmement, la mise à jour des ensembles *P-dominant* et *P-dominé* de chaque action dans le système d'information et enfin la mise à jour des approximations inférieures et supérieures des unions des classes de décision. L'algorithme proposé minimise le temps de calcul lorsqu'une action entre ou quitte le système d'information sans affecter la qualité des règles de décision inférées.

Afin d'inférer des règles de décision robustes, nous choisissons de généraliser l'algorithme présenté dans Li et al. (2013) afin de considérer l'entrée simultanée d'un ensemble d'actions.

4 MAI2P : Méthode de classification incrémentale

Cette section présente la méthode MAI2P que nous avons proposée pour la prédiction incrémentale et périodique de la classe de décision Cl_i à laquelle une action x est susceptible d'appartenir. Cette méthode se compose de trois phases (cf. Figure 1).

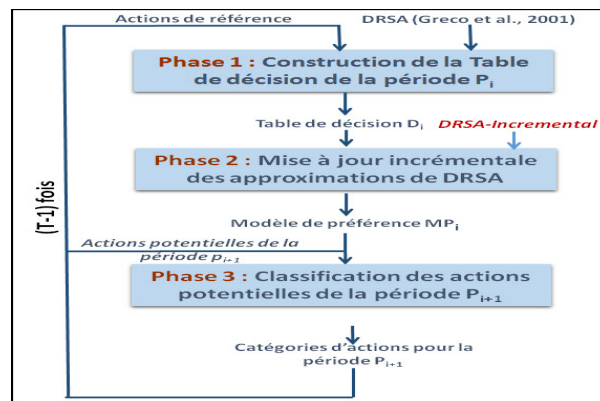


FIG. 1 – Description générale de la méthode MAI2P ($P_i =$ période i ; $T =$ nombre de périodes)

4.1 Phase 1 : Construction de la table de décision de la période P_i

Cette phase est composée de trois étapes :

4.1.1 Etape 1 : Construction d'un ensemble des "Actions de référence"

Cette étape consiste à définir un ensemble d'apprentissage contenant un nombre suffisant d'exemples représentatifs pour chacune des classes de décision prédéfinies. Afin de respecter la terminologie utilisée dans l'approche DRSA, nous appelons les exemples d'apprentissage, "Actions de référence". La construction de cet ensemble s'effectue par un ou plusieurs décideurs en fonction de leur expertise et leur expérience. D'un point de vue psychologique (Miller, 1956), un décideur humain se caractérise par une capacité cognitive représentant la limite supérieure à laquelle il peut associer ses réponses aux stimuli qui lui sont accordés. Ainsi, pour la construction de l'ensemble d'apprentissage, il est suffisant que les actions sélectionnées

soient représentatives et de qualité, quel que soit leur effectif. L'intervention des experts pour la construction de l'ensemble d'apprentissage permet d'obtenir des sous-ensembles équitables des "Actions de référence" et de surmonter le problème de données déséquilibrées.

La méthode *MAI2P* doit être appliquée sur les systèmes d'informations qui évoluent dans le temps, où l'ensemble des "Actions de référence" varie d'une période à une autre. Ainsi, chaque période P_i , le décideur doit définir un nouvel ensemble A'_i des "Actions de référence" qui se rajoute à l'ensemble des "Actions de référence", A_{i-1} , de toutes les périodes précédentes.

4.1.2 Etape 2 : Construction d'une famille cohérente de critères

Comparé à un attribut, un critère doit permettre de mesurer les préférences des décideurs selon un point de vue personnel (Mousseau et al., 1996). Dans ce travail, l'approche que nous adoptons est ascendante qui consiste à construire une famille de critères à partir d'une liste d'indicateurs susceptibles d'influencer l'opinion des décideurs concernant la caractérisation des actions. Ensuite, des réunions directes doivent être menées avec le décideur afin d'obtenir ses informations préférentielles sur chaque critère. Afin d'appliquer les points de vue préférentielles, nous adoptons une échelle qualitative.

4.1.3 Etape 3 : Classification de l'ensemble des "Actions de référence"

Cette étape consiste à la construction d'une table de décision D_i de la période P_i . C'est une matrice dont les colonnes représentent les " p " critères d'évaluation contenus dans F_i et dont les lignes forment un ensemble de " m " "Actions de référence" contenues dans A'_i . Le contenu de la matrice est la fonction d'évaluation $f_i(A_{j,i}, g_k)$ de chaque action $A_{j,i} \in A'_i$ sur chaque critère $g_k \in F_i$ tel que $i \in \{1..T\}$, $j \in \{1..m\}$ et $k \in \{1..p\}$. Les variables T , m et p sont respectivement le nombre de périodes à considérer pendant le processus de prédiction, la taille $|A'_i|$ de l'ensemble des "Actions de référence" définit à la $i^{\text{ème}}$ période et la taille de l'ensemble $|F_i|$ de la famille de critères. La dernière colonne de la table contient la décision d'affectation de chaque "Action de référence" dans l'une des N classes de décision.

4.2 Phase 2 : Mise à jour incrémentale des approximations de DRSA

Cette phase applique notre algorithme *DRSA-Incremental* (Bouzayane et Saad, 2017) sur la table de décision D_i construite pendant la période P_i afin d'en inférer un modèle de préférence, MP_i , susceptible de classer chaque action dans l'une des classes de décision prédéfinies.

Cette phase est appliquée dès que la table de décision est complète. Elle considère l'ensemble des "Actions de référence", A_{i-1} , de toutes les périodes précédentes et l'ensemble des "Actions de référence", A'_i , de la période P_i . L'algorithme *DRSA-Incremental* est déclenché dès l'insertion de l'ensemble A'_i dans la table de décision. Il est composé de quatre étapes :

1. Calculer les *unions supérieures* et *inférieures* de chacune des classes de décision Cl_{i-1} .
2. Calculer les *ensembles dominants* et *dominés* pour chaque action insérée $x^+ \in A'_i$.
3. Mettre à jour les *ensembles dominants* et *dominés* pour chaque action $A_{j,i-1} \in A_{i-1}$.
4. Mettre à jour les approximations de chacune des unions de classes de décision.

La sortie de la phase 2 est un modèle de préférence permettant de classer les "Actions potentielles" pendant la période P_{i+1} .

4.3 Phase 3 : Classification des “Actions potentielles” de la période P_{i+1}

La troisième phase exploite les règles de décision précédemment inférées afin d’attribuer chacune des “Actions potentielles” dans l’une des N classes de décision prédéfinies. Une “action potentielle” est une action susceptible d’être classée dans l’une des classes de décision.

Cette phase s’exécute pendant la période P_{i+1} tout au long du processus de prédiction tel que $i \in \{2, \dots, T\}$. Elle commence par l’évaluation de toutes les “Actions potentielles” sur l’ensemble de critères construits. Ensuite, il s’agit d’appliquer les règles de décision inférées pendant la période P_i afin de les affecter dans les classes de décision prédéfinies.

La méthode *MAI2P* s’exécute périodiquement tout au long du processus de prédiction : la première et la deuxième phases se déroulent pendant toutes les périodes P_i tel que $i \in \{1, \dots, T-1\}$ alors que la troisième se déroule pendant la période P_i ; tel que $i \in \{2, \dots, T\}$.

5 Expérimentation et évaluation de la méthode *MAI2P*

Nous avons traité le cas d’un MOOC (formation en ligne et gratuite) Français qui a duré 5 semaines et accédé par 2360 apprenants. L’objectif est la prédiction hebdomadaire de la classe de décision à laquelle appartiendra un apprenant : Cl_1 des “Apprenants en risque d’abandon”; Cl_2 des “Apprenants en difficulté” mais qui sont actifs; et Cl_3 des “Apprenants leaders”.

- *Phase1*. Nous avons construit, avec l’aide de l’équipe pédagogique, quatre ensembles des “Apprenants de référence” A'_i tel que $i \in \{1, 2, 3, 4\}$ et $|A'_i| = 30$. Ensuite, une famille cohérente de 11 critères a été définie dont 8 sont statiques (exp. niveau d’études) et 3 sont dynamiques (exp. le nombre hebdomadaire de messages). Enfin, à la fin de chaque semaine une table de décision est construite.
- *Phase2*. Cette phase a été appliquée à la fin de chaque semaine S_i une fois que la table de décision D_i est complète tel que $i \in \{1, 2, 3, 4\}$ en appliquant l’algorithme *DRSA-Incremental* pour la mise à jour incrémentale des règles de décision.
- *Phase3*. Cette phase était appliquée au début de chaque semaine S_i du MOOC en appliquant le modèle de préférence inféré à la fin de la semaine S_{i-1} pour la classification de l’ensemble d’apprenants potentiels tel que $i \in \{2, 3, 4, 5\}$.

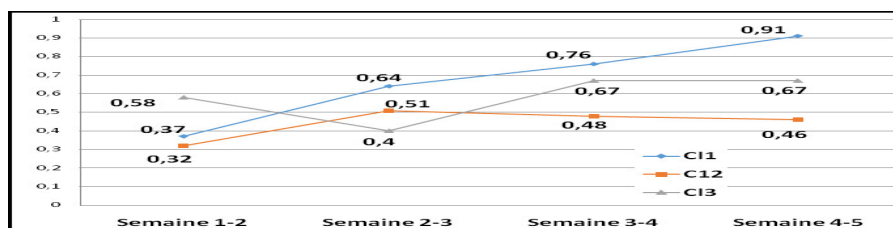


FIG. 2 – Qualité de la prédiction (*F-mesure*) durant les semaines du MOOC

La Figure 2 met l’accent sur la variation de la *F-mesure* des trois classes de décision Cl_1 , Cl_2 et Cl_3 d’une semaine à une autre.

- La *F-mesure* de la classe Cl_1 , des “Apprenants de risque”, augmente au cours du temps. En effet, le MOOC est connu par les *lurkers*. Ces apprenants restent actifs au bout de la première semaine mais en ayant une intention préalable d’abandonner la formation.

Ce type d'apprenants dévalorise la performance du modèle de prédiction qui est basé sur le profil et le comportement de l'apprenant et non pas sur son intention.

- La F-mesure de la classe Cl_3 , des "Apprenants leaders" augmente progressivement au cours du temps. En effet, d'une semaine à une autre, les apprenants multiplient leur participation au forum ce qui donne une information plus ample sur leur profils. Aussi, les évaluations proposées par le MOOC sont de plus en plus complexes d'une semaine à une autre ce qui permet d'une vision plus précise sur les compétences des apprenants.

6 Conclusion

Dans ce papier, nous avons proposé une méthode de classification multicritère et incrémentale *MAI2P* basée sur l'approche DRSA pour la prédiction périodique de la classe de décision à laquelle une action est susceptible d'appartenir. La méthode *MAI2P* est composée de trois phases : la construction d'une table de décision ; l'inférence d'un modèle de préférences en appliquant l'algorithme *DRSA-Incremental* et la prédiction de la classe de décision à laquelle appartiendra chaque action. Les expérimentations de la méthode *MAI2P* sur un MOOC Français ont démontré une qualité de prédiction satisfaisante qui atteint une F-mesure = 0.66.

Références

- Bouzayane, S. et I. Saad (2017). Incremental updating algorithm of the approximations in drsa to deal with the dynamic information systems of moocs. *In the international conference on Knowledge Management, Information and Knowledge Systems (KMIKS)*, 55–66.
- Greco, S., B. Matarazzo, et R. Slowinski (2001). Rough sets theory for multicriteria decision analysis. *EJOR* 129(1), 1–45.
- Greco, S., R. Slowinski, J. Stefanowski, et M. Zurawski (2004). Incremental versus nonincremental rule induction for multicriteria classification. *Transaction on Rough Sets II*, 33–53.
- Li, S., T. Li, et D. Liu (2013). Dynamic maintenance of approximations in dominance-based rough set approach under the variation of the object set. *Int. J. Intell. Syst.* 28, 729–751.
- Miller, G. A. (1956). The magical number seven, plus or minus two : some limits on our capacity for processing information. *Psychological review* 63(2), 81.
- Mousseau, B. R., VINCENT, et B. Roy (1996). A theoretical framework for analysing the notion of relative importance of criteria. *J. Multi-Criteria Decis. Anal* 5, 145–159.
- Roy, B. (1985). *Méthodologie multicritère d'aide à la décision*. Economica.

Summary

This paper proposes a method based on the rough set theory and dedicated to the incremental supervised learning in a context of unbalanced data. This method consists of three phases: the construction of a decision table, the inference of a set of decision rules, and the classification of each potential action in one of the predefined decision classes. The *MAI2P* method is validated in the context of MOOC (*Massive Open Online Course*).