

Projet RODM - Classification associative

(à effectuer seul ou en binôme, langage de programmation libre, Julia recommandé)

Nous considérons une liste de 32651 individus.

Comme illustré en Table 1, pour chaque individu, nous connaissons :

- 14 caractéristiques ;
- s'il gagne plus de 50k\$ ou non (classe de l'individu).

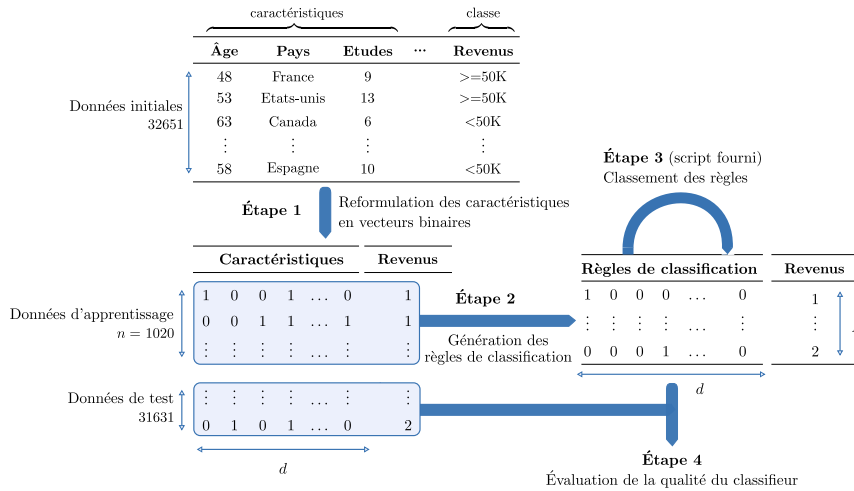
caractéristiques				classe
Âge	Pays	Etudes	...	Revenus
48	France	9	...	<50K
⋮	⋮	⋮	⋮	⋮
57	Espagne	10	...	>=50K

TABLE 1 – *Extrait des données initiales. Chaque ligne correspond à un individu. La dernière colonne contient la valeur ckd, s'il est malade.*

L'objectif de ce projet est de mettre en œuvre la technique de classification associative de Bertsimas, Chang et Rudin [1], présentée en cours, afin de diagnostiquer si les individus sont atteints d'une maladie rénale.

Le travail minimum demandé pour ce projet est :

- la réalisation des 4 étapes représentées dans le schéma ci-dessous ;
- le traitement d'au moins une des questions d'ouverture de l'étape 5 (totalement ou partiellement suivant la difficulté).



Le code fourni avec ce sujet facilite l'implantation en Julia. Il contient deux fichiers :

- **main.jl** : contenant l'inclusion des packages et l'appels aux principales fonctions (binarisation, création des règles, tri des règles et calculs des performances) ;
- **function.jl** : contenant le code des différentes fonctions dont une partie est à compléter.

Afin d'exécuter le programme il faut taper se déplacer dans le répertoire **src** et taper la commande suivante :

```
julia main.jl
```

Étape 1 Reformulation des caractéristiques sous la forme de vecteurs binaires (méthode createFeatures)

La méthode de classification considérée nécessite que chaque individu soit représenté par un vecteur de caractéristiques binaires. Cependant, la plupart des caractéristiques ne le sont pas.

1. Proposer une représentation binaire pertinente des caractéristiques.

Remarque : Vous n’êtes pas obligé d’utiliser toutes les caractéristiques et plus vous aurez de caractéristiques, plus la résolution des problèmes d’optimisation risque d’être longue ou non optimale.

2. Reformuler les données initiales, situées dans le fichier “data/adult.csv”, en utilisant cette représentation. Pour cela vous pourrez vous aider de la méthode `createColumns` qui permet de binariser une colonne numérique à partir d’intervalles.
3. Créer le fichier de données “data/adult_train.csv” contenant $\frac{1}{32}$ des individus binarisés, choisi aléatoirement, ainsi qu’un fichier “data/adult_test.csv” contenant le reste des individus binarisés.

Étape 2 Génération des règles de classification (méthode `createRules`)

Implanter l’algorithme vu en cours pour obtenir l’ensemble des règles de classification dans le fichier “res/adult_rules.csv”. Pour cela, vous pourrez compléter la méthode `createRules` du fichier `functions.jl` où les valeurs des différents paramètres de la méthode sont fixés.

Étape 3 Classement des règles (méthode `sortRules`)

Utiliser la méthode `sortRules` du fichier `functions.jl` afin de déterminer un classement optimal des règles. Après exécution, les règles ordonnées se trouveront dans le fichier “res/adult_ordered_rules.csv”.

Remarque : Par défaut le temps de résolution est limité à 300 secondes mais vous pouvez l’augmenter afin d’améliorer les performances (ou si aucune solution entière n’est trouvée).

Étape 4 Évaluation de la qualité du classifieur (méthode `showStatistics`)

Calculer la précision et le rappel du classifieur sur les données d’apprentissage et de test. Ces résultats vous semblent-ils pertinents ?

Étape 5 Questions d’ouverture

1. Comparer les performances du classifieur lorsque d’autres vecteurs de caractéristiques sont considérés.
2. Appliquer la méthode à un autre corpus (par exemple un de ceux figurant au [lien suivant](#)). Le corpus que vous choisirez ne doit pas avoir été considéré les années précédentes (*i.e.*, différent de titanic, haberman et kidney) ;
3. L’étape 2 consiste à générer les règles qui maximisent le support tout en minimisant la couverture. L’algorithme `RuleGen` fait le choix de fixer une borne supérieure sur la couverture afin de ne pas se trouver face à un problème bi-objectif. Dans cette question nous allons résoudre directement le problème bi-objectif. Pour ce faire, utiliser les packages `vOptGeneric` et `MultiJuMP` de Julia qui permettent de résoudre facilement des problèmes de ce type en retournant une liste de solutions (chaque solution correspondant, ici, à une règle de classification). Étudier l’intérêt de cette approche par rapport à `RuleGen` (nombre de règles obtenues, temps de calcul, performance du classifieur, ...).
4. Afin d’accélérer la résolution, résoudre uniquement les relaxation linéaires des problèmes d’optimisation de l’étape 2 et/ou 3 (*i.e.*, résolution de problèmes continus au lieu de problèmes entiers). Vous arrondirez les solutions obtenues et comparerez les performances. Évaluer ensuite dans quelle mesure cela vous permet d’augmenter le nombre de caractéristiques considérées ou de règles générées, puis comparez une nouvelle fois les performances.
5. Certaines des inégalités du problème de l’étape 3 ont pour principal objectif de renforcer la relaxation linéaire de la formulation. Déterminer celles qui sont les plus utiles (si elles s’avèrent effectivement utiles).
6. Évaluer l’impact de la limite de temps de résolution de l’étape 3 sur le gap et sur la qualité du classifieur obtenu.
7. Tout autre contribution qui peut vous sembler pertinente du point de vue des performances ou de l’application peut également faire l’objet d’une question d’ouverture.

Étape 6 Rendu

A l'issue de ce projet (data limite de rendu le 14/03), vous rendrez vos fichiers sources, les fichiers de données générés ainsi qu'un rapport qui contiendra notamment :

- Une description argumentée de la représentation binaire que vous avez choisie ;
- Le nombre de règles générées, le temps de calcul associé, le solveur choisi, ainsi que sa version, et les caractéristiques de la machine utilisée (type et nombre de processeurs et mémoire RAM) ;
- Le nombre de règles du classifieur et le temps de calcul pour les ordonner ;
- Les performances de votre classifieur (précision et rappel).
- Un tableau représentant les règles ordonnées de votre classifieurs ainsi que leur classe accompagnés d'une description de la façon dont le classifieur fait ses choix.
- Le travail effectué sur les questions d'ouvertures.
- La liste des opérations à effectuer pour reproduire vos expériences (quels programmes exécuter, avec quels paramètres, dans quel ordre, ...).

Références

- [1] Allison Chang, Dimitris Bertsimas, and Cynthia Rudin. An integer optimization approach to associative classification. In *Advances in neural information processing systems*, pages 269–277, 2012.