

Classification associative

Alhousseynou BALL* and Bakary SIDIBE*

*Etudiants en M2 Datascience - Ecole Polytechnique X

11th March 2020

Contents

1	Données et représentation binaire	3
1.1	Présentation des données	3
1.2	Représentation binaire	3
2	Généralisation des règles et ordonnancement	6
3	Performances	7
4	Questions d'ouvertures	8
4.1	performances du classifieur sur d'autres variables	8
4.2	Application sur d'autres données	9
5	Étapes pour reproduire ce travail	10

Introduction

Un des enjeux majeurs des scientifiques de données est de mettre en place de puissants algorithmes afin de synthétiser les informations pour optimiser la prise de décisions. C'est dans ce cadre que la classification associative a été mise en place. Elle est une sorte de modèle d'apprentissage supervisé qui utilise des règles d'association pour attribuer une valeur cible. L'article de Bertsimas et al. (2012) propose une nouvelle approche de Mixed Integer Optimization (MIO) appelée Ordered Rules for Classification pour cette tâche de classification. Cette méthode comporte deux parties:

- La partie explore une frontière particulière des solutions dans l'espace des règles
- La deuxième partie apprend à établir un classement optimal pour les règles permettant ainsi de construire une liste de décision qui est simple et perspicace.

L'objectif de notre document est de mettre en pratique cet algorithme sur des données réelles et puis évaluer ses performances.

Ce document est structuré en cinq principales parties. La première partie décrit brièvement la base de données et propose une représentation binaire des variables à retenir pour la prédiction de la variable cible. La deuxième partie porte sur la généralisation des règles de classification et l'ordonnancement de ces dernières. La troisième partie évalue les performances obtenues. La quatrième partie traitera des questions d'ouverture. Enfin, la cinquième partie donnera les étapes à faire afin de reproduire ce travail.

1 Données et représentation binaire

1.1 Présentation des données

Les données sont composées de 189 patients diagnostiqués ou non de la maladie rénale chronique(ckd or notckd). On cherche à prédire si une personne est atteinte de la maladie rénale chronique ou pas en se basant sur certaines variables.

variables	description	variables	description	variables	description
age	age	bgr	blood glucose random	rbcc	red blood cell count
bp	blood pressure	bu	blood urea	htn	hypertension
sg	specific gravity	sc	serum creatinine	dm	diabetes mellitus
al	albumin	sod	sodium	cad	coronary artery disease
su	sugar	pot	potassium	appet	appetite
pc	pus cell	hemo	hemoglobin	pe	pedal edema
pcc	pus cell clumps	pcv	packed cell volume	ane	anemia
ba	bacteria	wbcc	white blood cell count	class	class

(a)
(b)
(c)

Figure 1: Description Kidney data

	age	bp	sg	al	su	pc	...	appet	pe	ane	class
0	48	70	1.005	4	0	abnormal	...	poor	yes	yes	ckd
1	53	90	1.020	2	0	abnormal	...	poor	no	yes	ckd
2	63	70	1.010	3	0	abnormal	...	poor	yes	no	ckd
3	68	70	1.015	3	1	normal	...	poor	yes	no	ckd
4	68	80	1.010	3	2	abnormal	...	poor	yes	no	ckd

Figure 2: Les 5 premières observation

1.2 Représentation binaire

Avant d'appliquer un modèle de classification, il faut dans un premier temps choisir les variables pertinentes puis dans un second temps rendre les variables sélectionnées binaires. Pour le choix des variables pertinentes, une étude de corrélation entre la variable cible et les prédicateurs a été menée. Étant donné que les variables sont mixtes, deux mesures de corrélation sont utilisées: rapport de corrélation et la corrélation de Pearson. Une fois cette première sélection faite, nous allons étudier la multicollinearité entre les variables avec l'algorithme VIF pour sélectionner finalement les variables qui garantissent une bonne performance du modèle.

- rapport de corrélation: Le rapport de corrélation mesure le niveau de liaison entre une variable qualitative et une variable quantitative. Il correspond au rapport entre la variance entre les groupes (groupe de positifs et groupe de négatifs) et la variance totale. Sa valeur varie entre 0 et 1. Un rapport de corrélation proche de 1 révèle une forte corrélation entre les variables en jeu, et un rapport de corrélation proche de 0 révèle une faible liaison entre les variables. Le choix des variables à retenir se base sur le test de nullité du rapport de corrélation. Il suppose que la quantité $F = \frac{r(n-p)}{(p-1)(1-r)}$ suit une loi de Fischer de paramètres $p-1$ et $n-p$ avec r le rapport de corrélation, n le nombre total d'individus et p le nombre de groupes. Si F est supérieur au quantile d'ordre 0,99 de la loi de Fisher à $p-1$ et $n-p$ degrés de liberté, alors l'hypothèse nulle de nullité du rapport est rejetée au seuil de 1% et donc on conclut en la corrélation entre la variable à expliquer et la variable explicative. On a $n=189$ et $p=2$, alors le quantile d'ordre 0,99 vaut $F_{1,187}^{0,99} = 6.78$. En se référant sur la figure 3, la variable **pot** ne se sera pas considérée dans la modélisation.

Variables	Fisher	Variables	Fisher
hemo	363.49	bu	81.68
pcv	362.36	bgr	67.14
sg	272.41	su	47.80
al	253.46	age	31.66
rbcc	189.74	wbcc	26.91
sc	98.07	bp	26.65
sod	85.78	pot	2.40

(a)
(b)

Figure 3: Fisher corrélation

- corrélation de Pearson (Khi2): La corrélation de Pearson permet d'étudier la liaison entre deux variables qualitatives. Étant donné un tableau représentant la distribution jointe de deux variables, la statistique de test est donnée par : $\chi^2 = n(\sum_{i,j} \frac{n_{i,j}^2}{n_{i.}n_{.j}} - 1)$ avec n_{ij} est l'effectif contenu dans case repérée par la ligne i et la colonne j , $n_{i.}$ est l'effectif marginal de la ligne i , $n_{.j}$ est celui de la colonne j , et n est l'effectif total. On remarque que, d'après la figure 4, l'ensemble des variables qualitatives sont corrélées avec la variable cible.

variables	stats	pvalue	criticals
htn	120.07	0.0	6.63
dm	96.59	0.0	6.63
pc	63.66	0.0	6.63
appet	46.02	0.0	6.63
pe	43.94	0.0	6.63
ane	33.90	0.0	6.63
pcc	28.16	0.0	6.63
cad	24.45	0.0	6.63
ba	19.04	0.0	6.63

Figure 4: Corrélation Khi2

- Variance inflation factor: Le VIF quantifie le niveau de multicolinéarité entre les variables à partir de la valeur du R^2 simple obtenue après régression multiple des variables explicatives sur la variable à expliquer. Il s'exprime par la relation : $VIF = \frac{1}{1-R^2}$ et la multicolinéarité est forte lorsque $VIF > 10$. En se référant à la figure 5, on n'a pas de problème de multicolinéarité forte. Toutes les variables restantes peuvent être utilisées pour la classification. Cependant, en nous basant sur la précision et le recall, nous avons choisi les variables qui minimisent l'erreur de classement. Ces variables sont: hemo, pcv, sg, htn, dm, pc.

variables	VIF	variables	VIF	variables	VIF
hemo	5.7880	dm	3.4578	pcc	1.9609
pcv	5.6128	rbcc	3.1112	sod	1.8676
sc	4.7766	pc	2.8822	ba	1.7982
htn	4.6965	sg	2.6139	cad	1.5307
bu	4.4951	ane	2.5461	bp	1.3363
al	4.0704	pe	2.5270	wbcc	1.3101
bgr	3.6284	appet	2.0966	age	1.3092
su	3.4614				

(a)
(b)
(c)

Figure 5: Variance inflation factor

2 Généralisation des règles et ordonnancement

Dans un premier temps, il était question d'implémenter l'algorithme du RuleGen qui fournit les règles de classification. Le nombre total de règles obtenu est de 79.

Algorithme RuleGen

Entrées mincov, iter_lim
pour toute classe y faire

- $(\mathcal{R}_Y, \bar{s}, \text{iter}, c_{\max}) \leftarrow (\emptyset, 0, 1, n)$
- Répéter
 - si iter = 1 alors
 - $\bar{s}, b \leftarrow \text{Résoudre } P(y, c_{\max})$
 - iter = iter + 1
 - $\mathcal{R}_Y \leftarrow \mathcal{R}_Y \cup b$
 - Ajouter la contrainte(*)
 - si iter < iter_lim alors
 - stemp, $b \leftarrow \text{Résoudre } P(y, c_{\max})$
 - si sTemp < \bar{s} alors
 - $c_{\max} \leftarrow \min(c_{\max} - 1, \sum_{i=1}^n x_i)$
 - iter $\leftarrow 1$
 - sinon
 - iter $\leftarrow \text{iter} + 1$
 - sinon
 - $c_{\max} \leftarrow c_{\max} - 1$
 - iter $\leftarrow 1$
 - jusqu'à Cmax < n mincov
- retourner \mathcal{R}_Y

Le nombre de règles du classifieur est donné par la figure 6.

x1	x2	x3	x4	x5	x6	x7	x8	x9	x10	x11	x12	x13	x14	x15	x16
1	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0
1	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0
1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Figure 6: Les règles du classifieur: kidney data

Interprétation :

- si $dm = \text{'yes'}$ et $hemo \in [11.3-13.9]$, on prédit que l'individu est malade('ckd')
- sinon, si $htn = \text{'yes'}$ et $hemo \in [11.3-13.9]$, on prédit que l'individu est malade('ckd')
- sinon, si $sg \in [0-1.015]$, on prédit que l'individu est malade('ckd')
- sinon, si $hemo \in [0-11.3]$, on prédit que l'individu est malade('ckd')
- sinon, si $sg \in [1.015-1.02]$, on prédit que l'individu est malade('ckd')
- sinon, on prédit que l'individu n'est pas malade('notckd')

Les caractéristiques de la machine sont les suivants et les temps de calculs sont les suivants:

Système	
Processeur :	Intel(R) Core(TM) i7-2620M CPU @ 2.70GHz 2.70 GHz
Mémoire installée (RAM) :	8,00 Go (7,88 Go utilisable)
Type du système :	Système d'exploitation 64 bits, processeur x64

Figure 7: informations sur la machine

Root node processing (before b&c):
Real time = 0.94 sec. (98.24 ticks)
Parallel b&c, 4 threads:
Real time = 0.00 sec. (0.00 ticks)
Sync time (average) = 0.00 sec.
Wait time (average) = 0.00 sec.
Total (root+branch&cut) = 0.94 sec. (98.24 ticks)

(a) Generating the rules

CPXPARAM_TimeLimit	300
Root node processing (before b&c):	
Real time = 22.64 sec. (8192.40 ticks)	
Parallel b&c, 4 threads:	
Real time = 277.44 sec. (105214.02 ticks)	
Sync time (average) = 29.98 sec.	
Wait time (average) = 0.00 sec.	
Total (root+branch&cut) = 300.08 sec. (113406.42 ticks)	

(b) Sorting the rules

Figure 8: Performance sur Kidney data

Le solveur utilisé est le CPLEX version 12.9.

3 Performances

Les performances du classifieur seront évaluées à partir de deux métriques: précision et recall.

Le classifieur prédit bien la classe des individus, car tous les individus étiquetés comme étant malades sont réellement malades.

Class	Prec.	Recall	Size
0	1.0	1.0	77
1	1.0	1.0	49
avg	1.0	1.0	
w. avg	1.0	1.0	

(a) Train results

Class	Prec.	Recall	Size
0	1.0	1.0	38
1	1.0	1.0	25
avg	1.0	1.0	
w. avg	1.0	1.0	

(b) Test results

Figure 9: Performance sur Kidney data

4 Questions d'ouvertures

4.1 performances du classifieur sur d'autres variables

Le classifieur donne de bons résultats avec les nouvelles variables. Cependant, comparé avec les résultats obtenus ci-dessus, il semble être moins performant lorsqu'on utilise d'autres vecteurs de caractéristiques. La précision et le recall ont diminué globalement de 2%.

Class	Prec.	Recall	Size	Class	Prec.	Recall	Size
0	0.98	1.0	79	0	0.97	1.0	36
1	1.0	0.96	47	1	1.0	0.96	27
avg	0.99	0.98		avg	0.99	0.98	
w. avg	0.98	0.98		w. avg	0.98	0.98	

(a) Train results

(b) Test results

Figure 10: Performance sur Kidney data

Le nombre de règles du classifieur sont devenus plus importants(12 règles contre 6).

x1	x2	x3	x4	x5	x6	x7	x8	x9	x10	x11	x12	x13
1	0	0	0	0	0	0	1	0	0	0	0	0
0	0	0	0	0	1	0	0	0	1	0	0	0
1	0	0	0	0	1	0	0	0	0	0	0	1
0	1	0	0	0	1	0	0	0	0	0	0	0
0	1	0	0	1	0	0	0	0	1	0	0	0
1	0	0	0	0	0	1	0	0	0	0	1	0
1	0	1	0	0	0	1	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	1	1
1	0	0	0	0	0	0	0	0	0	0	1	0
0	1	0	1	0	0	0	0	0	1	0	0	0
1	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0

Figure 11: Les règles du classifieur: kidney data

4.2 Application sur d'autres données

La méthode décrite ci-dessus a été appliquée sur les données de vote du Congrès américain de 1984. La description des variables est fournie en annexe à la figure 13. La base de données disponible sur [UCI Machine Learning](#) comporte des valeurs manquantes. Ainsi, une imputation par le mode a été effectuée.

Le classifieur utilise les 6 règles de la figure 12a.

x1	x2	x3	x4	x5	x6	Class	Prec.	Recall	Size	Class	Prec.	Recall	Size
1	1	0	1	1	0	0	0.98	0.96	180	0	0.97	0.97	87
1	1	0	1	0	1								
1	1	0	0	0	1	1	0.93	0.96	110	1	0.95	0.95	58
0	0	1	0	0	0	avg	0.95	0.96		avg	0.96	0.96	
1	1	0	0	0	0								
0	0	0	0	0	0	w. avg	0.96	0.96		w. avg	0.96	0.96	

(a) règles du classifieur
(b) Train results
(c) Test results

Figure 12: Performance sur house votes

Interprétation

- si physician = 'yes', salvador = 'yes' et education = 'yes', on prédit que l'individu votera pour les républicains
- sinon, si physician = 'yes', salvador = 'yes' et nicaraguan = 'yes', on prédit que l'individu votera pour les républicains
- sinon, si physician = 'yes' et nicaraguan = 'yes', on prédit que l'individu votera pour les républicains
- sinon, si adoption = 'yes', on prédit que l'individu votera pour les démocrates
- sinon, si physician = 'yes', on prédit que l'individu votera pour les républicains
- sinon, on prédit que l'individu votera pour les démocrates

5 Étapes pour reproduire ce travail

Pour reproduire le travail, voici les étapes à suivre:

Installation

1. installer:
 - (a) Julia
 - (b) CPLEX (v12.8 ou v12.9): attention, la 12.10 n'est pas compatible avec Julia pour l'instant.
2. installer les packages: tapez "julia" dans une console, puis:
 - (a) `using Pkg`
 - (b) `Pkg.add("JuMP")`
 - (c) `Pkg.add("DataFrames")`
 - (d) `Pkg.add("CSV")`
 - (e) `ENV["CPLEX_STUDIO_BINARIES"] = "votre chemin d'installation de cplex (plus d'infos en dessous)"`
 - (f) `Pkg.add("CPLEX")`

Le chemin à indiquer pour CPLEX :

Linux

```
ENV["CPLEX_STUDIO_BINARIES"] = "/path/to/cplex/bin/x86-64_linux"
```

OSX

```
ENV["CPLEX_STUDIO_BINARIES"] = "/path/to/cplex/bin/x86-64_osx"
```

Windows

```
ENV["CPLEX_STUDIO_BINARIES"] = "C:/IBM/CPLEX_Studio128/cplex/bin/x64_win64"
```

““

Kidney data

1. Utiliser les variables **class**, **htn**, **dm**, **hemo**, **pcv**, **sg**, **pc** puis les rendre binaire en respectant ces intervalles(quartiles):
 - (a) hemo: [0.0-11.3[, [11.3-13.9[, [13.9-15.5[, [15.5-Inf[
 - (b) pcv: [0.0-34.0[, [34.0-42.0[, [42.0-48.0[, [48.0-Inf[
 - (c) sg: [0.0-1.015[, [1.015-1.02[, [1.02-1.025[, [1.025-Inf[
2. ouvrir **julia** puis se placer dans le dossier **src** et exécuter la commande **include('main.jl')**

Pour reproduire les questions d'ouvertures, appliquer la procédure décrite ci-dessus en utilisant ces vecteurs caractéristiques : **al**, **sc**, **rbcc**, **htn**, **dm** et la binarisation suivante:

- al: [0.0-1.0[, [1.0-Inf[
- sc: [0.0-0.8[, [0.8-1.1[, [1.1-2.2[, [2.2-Inf[
- rbcc: [0.0-4.0[, [4.0-4.8[, [4.8-5.5[, [5.5-Inf[

housevotes data

Les variables utilisées sont **physician,adoption,salvador,education,nicaraguan**

Conclusion

Dans ce projet, nous avons mis en pratique la méthode Ordered Rules for Classification (ORC) proposée par l'article de Bertsimas et al. (2012). L'avantage majeur de cette méthode par rapport à d'autres méthodes de classification en machine learning est son interprétabilité. De plus, les performances sont globalement très bonnes.

Annexes

variables	description	variables	description
class	Class	missile	mx-missile
hand	handicapped-infants	immigration	immigration
water	water-project-cost-sharing	synfuels	synfuels-corporation-cutback
adoption	adoption-of-the-budget-resolution	education	education-spending
physician	physician-fee-freeze	superfund	superfund-right-to-sue
salvador	el-salvador-aid	crime	crime
religious	religious-groups-in-schools	exports	duty-free-exports
satellite	anti-satellite-test-ban	administration	export-administration-act-south-africa
nicaraguan	aid-to-nicaraguan-contras		

(a)
(b)

Figure 13: Description housevotes data

References

Bertsimas, D., A. Chang, and C. Rudin (2012). An Integer Optimization Approach to Associative Classification.