EXPOSE DE FOUILLE DE DONNEES ET RECHERCHE D'INFORMATION,

THÈME: MÉTHODE D'APPRENTISSAGE SUPERVISÉ: PERFORMANCE DES ÉTUDIANTS

> Réalisé par : ADOUM OKIM BOKA Encadreur : Mme Nguyen THI MINH HUYEN

> > ntmhuyen@gmail.com



August 20, 2019

SOMMAIRE



INTRODUCTION

PRÉSENTATION DES DONNEES

Présentation de données

ANALYSE EXPLORATOIRE

Methode Factorielle
Méthode de classification

CHOIX D'UNE MÉTHODE D'APPRENTISSAGE SUPERVISE

Méthodes

PRÉPARATION DES DONNÉES

Pré traitement des données

CONSTRUCTION ET ÉVALUATION DES MODÈLES

Matrice de confusion

Indicateurs de performances

Description de l'arbre

Performances de la solution θ = 0.7

COMPARAISON AVEC UNE AUTRE MÉTHODE

CONCLUSION

EXPOSE DE FOUILLE DE DONNEES ET RECHERCHE D'INFORMATION, THÈME : MÉT

INTRODUCTION



L'analyse des données est une famille de méthodes statistiques dont les principales caractéristiques sont d'être multidimensionnelles, descriptives et prédicatives.

PRÉSENTATION DES DONNEES



Présentation de données

 Les	attributs de doni	nees comprennent	les notes d	ies eleves;
Les	caractéristiques	démographiques.	sociales et	scolaires ont

- été collectés -> rapports et de questionnaires;
- ☐ Performance dans deux matières distinctes (Por Maths).

Résumé : Prédisez les performances des élèves dans l'enseignement secondaire (lycée).

Caractéristiques du jeu de données:	Multivarié	Nombre d'instances:	649	Surface:	Social
Caractéristiques d'attribut:	Entier	Nombre d'attributs:	33	Date du don	2014-11-27
Tâches associées:	Classification, régression	Valeurs manquantes?	N/A	Nombre de visites sur le Web:	425362

ANALYSE EXPLORATOIRE



Cette section est divisé en deux parties:

- ☐ Réduction du nombre de variables (dimension) : méthodes factorielles:
- ☐ Réduction du nombre d'individus en regroupant par classes : méthodes de classification.

ANALYSE EXPLORATOIRE

Methode Factorielle



Methode Factorielle:

Les données sont classifiés en deux types qui sont:

- les valeurs discrètes dont nous allons appliqués l'analyse de correspondance multiples;
- ☐ les valeurs continues pour l'analyse en composante principal.

ANALYSE EXPLORATOIRE

Méthode de classification



Méthode de classification:

La classification (clustering) est une méthode mathématique d'analyse de données qui permet de faciliter l'étude d'une population d'effectif important (personnes, animaux, plantes, malades, gènes,...), dans l'optique de les regroupe en plusieurs classes de telle sorte que les individus d'une même classe soient le plus semblables possible et que les classes soient le plus distinctes possibles.

CHOIX D'UNE MÉTHODE D'APPRENTISSAGE SUR VISE

Analyse exploratoire
☐ méthodes factorielles;
☐ méthodes de classification.
L'attribut cible G3 a une forte corrélation avec les attributs G2 et G1.
☐ Méthode de Bayes naïf;
☐ K plus proches voisins;
☐ Arbre de décision;
☐ Réseaux de neurones;
□ algorithme K-NN.

CHOIX D'UNE MÉTHODE D'APPRENTISSA SUPERVISE Méthodes

Méthodes

- ☐ Arbre de décision;
- □ CART(Classification And Régression Trees);
- □ algorithme K-NN;

PRÉPARATION DES DONNÉES



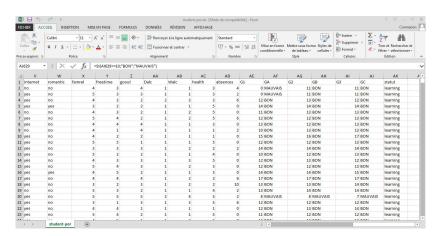
Objectif est de prédire la performance des élèves en s'appuyant principalement sur l'attribut G3. **pour ce faire, nous allons** transformer les attributs **G1**, **G2**, **G3** en des valeurs discrètes. **Afin quelles soit facilement segmentale**

PRÉPARATION DES DONNÉES

Pré-traitement des données



Pré traitement des données:



EXPOSE DE FOUILLE DE DONNEES ET RECHERCHE D'INFORMATION, THÈME : MÉT

CONSTRUCTION ET ÉVALUATION DES MODÈLES



Arbre de décision – La méthode CART

	ault title			O.	screte select examples 1			
□ III Dataset (student-por.x)	ls)		Discrete select examples 1 Parameters					
Discrete select exa		Attribute selection : statut Value selection : learning						
- Dupervised C		Results						
		520 selected exam	520 selected examples from 649					
								-
			Computation time : 0 ms. Created at 76/06/2019 18:57:50					
			M COME OF 2017 1012 (30)					
				Components				
Data visualization	Statistics	Nonparametric statistics	Instance selecti	Components on Feature construction	Feature selection	Regression	Factorial analysis	
Data visualization PLS	Statistics Clustering	Nonparametric statistics Spv learning	Instance selecti Meta-spv learnin	on Feature construction	Feature selection Scoring	Regression Association	Factorial analysis	
PLS	Clustering			on Feature construction	Scoring		Factorial analysis	
PLS		Spv learning		on Feature construction ng Spv learning assessment	Scoring		Factorial analysis	
PLS PLS Binary logistic regression BVM C4.5	Clustering As CS-CRT A CS-MC4 CS-VC	Spy learning A ID3 K-NN Linear dis	Meta-spv learnin	on Feature construction Spv learning assessment Final Multinomial Logistic Regression Final Maire bayes \$\frac{1}{40}\$ Naive bayes continuous	Scoring Prototype-NN Radial basis function R. Rnd Tree		Factorial analysis	
	Clustering As CS-CRT As CS-MC4	Spy learning A ID3 K-NN	Meta-spv learnin	on Feature construction g Spy learning assessment Multinomial Logistic Regression multinomial sayes	Scoring Prototype-NN Radial basis function		Factorial analysis	

CONSTRUCTION ET ÉVALUATION DES MÈLES Matrice de confusion

Matrice de confusion:

La matrice de confusion confronte les vraies valeurs et les valeurs prédites de GC sur les 520 observations ayant participé à l'apprentissage (growing + pruning). Elle est accompagnée du taux d'erreur qui est de 0.0577 dans notre exemple.

CONSTRUCTION ET ÉVALUATION DES MEÈLES

Indicateurs de performances

Indicateurs de performances:



EXPOSE DE FOUILLE DE DONNEES ET RECHERCHE D'INFORMATION, THÈME : MÉT

Description de l'arbre:

Tree description

Number of nodes	7
Number of leaves	4

Decision tree

- GB in [BON] then GC = BON (99,30 % of 286 examples)
- GB in [MAUVAIS]
 - famrel < 4,5000
 - GA in [MAUVAIS] then GC = MAUVAIS (61,54 % of 26 examples)
 - GA in [BON] then GC = BON (77,78 % of 18 examples)
 - famrel >= 4,5000 then GC = MAUVAIS (94,44 % of 18 examples)

Computation time: 63 ms.

CONSTRUCTION ET ÉVALUATION DES M

ÈLES
Performances de la solution $\theta = 0.7$

Performances de la solution θ = 0.7

Nous obtenons un taux d'erreur de 0.1163, avec un intervalle de confiance à 88% égal à [0.1312; 0.1398]. En comparant les différents taux d'erreur du TEST par rapport à (0, 0.7 et 1), nous constatons que l'arbre à 4 feuilles suffit largement pour assurer un niveau de performances satisfaisant

COMPARAISON AVEC UNE AUTRE MÉTHODE



Algorithme K-NN

Classifier performances CART

Error rate			0,0577				
Valu	ies pred	diction	Confusion matrix				
Value	Recall	1-Precision		BON	MAUVAIS	Sum	
BON	0,9544	0,0200	BON	440	21	461	
MAUVAIS	0,8475	0,2958	MAUVAIS	9	50	59	
			Sum	449	71	520	

Class	ifier	perfo	rmano	es	K-	NN
	Error ra	ite		0,0	0635	
Valu	ies pred	diction	Confusion matrix			
Value	Recall	1-Precision		BON	MAUVAIS	Sum
BON	0,9913	0,0597	BON	457	4	461
MAUVAIS	0,5085	0,1176	MAUVAIS	29	30	59
			Comm	407	24	E20

CONCLUSION



- Choix de la méthode;
- Stratégie pour obtenir un bon résultat;
- Nos résultats.

Merci pour votre aimable attention!

