

Rapport Projet

AGBLODOE Komi/Azat ALEKSANYAN

2/9/2020

Contents

1. Introduction	2
2 Présentation de l'application	3
2.1 Présentation succincte de Shiny	3
2.2 Rôles de l'application	3
2.3 Contraintes de l'application	3
2.4 Présentation générale	3
3 Algorithmes de Machine Learning	6
kNN /k-ppv/	6
CART - Classification et Regression Trees	6
rfRanger ???	7
Random Forest	7
XGBoost	7
glmNetL2 - Ridge Regression	7
glmNetL1 - Lasso Regression	8
AdaBoost Classification Trees ???	8
AdaBoost M1	8
Optimisation du modèle - Validation Croisée	9

1. Introduction

Dans le cadre de notre préparation au diplôme d'ingénieur en statistiques et sciences de données, nous avons été amenés à effectuer notre projet de cinquième année en partenariat avec l'entreprise BD (Becton et Dickinson) et sous l'encadrement de Didier Morel.

Il concerne le développement d'une application RShiny permettant de comparer et de tester des algorithmes de Machine Learning. Dans la mesure où il n'existe apparemment aucune application en ligne permettant de réaliser ce type de travail, il semble très intéressant de mettre en place un tel outil. Ainsi, les expérimentés comme les personnes sans trop de connaissances en apprentissage automatique pourront se servir de l'application pour avoir un aperçu de l'efficacité de différentes méthodes.

Dans ce projet, après avoir identifié le contexte, les besoins et les contraintes, nous avons procédé à la construction de notre application en créant deux fichiers dont un pour le serveur et le second pour l'interface utilisateur. Nous avons implémenté dans notre application, des algorithmes de classification et de régression comme AdaBoost Classification Trees, AdaBoost.M1, CART, les k plus proches voisins (kNN), le randomForest, RfRanger, XGBoost, glmNetL1, glmNetL2. Pour ces différents algorithmes, l'utilisateur a la possibilité de choisir les paramètres voulus. Nous avons ensuite procédé aux tests sur des jeux de données afin de comparer la performance des différents algorithmes de classification. Il faut noter aussi que notre application est capable de charger les jeux de données de type CSV et TSV.

2 Présentation de l'application

Notre application web est constituée d'un menu permanent qui permet de naviguer. Sur le site, une page d'accueil présente l'outil général à l'internaute.

2.1 Présentation succincte de Shiny

Shiny, est un package R qui facilite la construction d'applications web interactives depuis R permettant une sortie en html. Une application shiny se structure en deux parties à savoir :

-Un côté UI (User Interface) qui regroupe tous les éléments de mise en forme et d'affichage de l'interface utilisateur (affichage des inputs et des outputs).

-Un côté Serveur où sont exécutés les codes R qui servent à produire les éléments de sorties (graphiques, tables, traitements, etc..) et à les mettre à jour en cas de changement dans les valeurs d'entrées (inputs).

Les utilisateurs peuvent simplement manipuler l'application pour exécuter et afficher des résultats fournis par du code R. Notre application a été structurée en deux parties comme décrites ci-haut.

2.2 Rôles de l'application

Le travail qui nous a été demandé par Didier Morel est le développement d'une application avec Shiny. Le développement de cette application passe par l'implémentation des algorithmes de machine learning. Cette application doit être capable de:

- Tester l'efficacité des algorithmes (vitesse d'exécution)
- Tester la qualité des modèles fournis par les algorithmes (erreurs/taux de bonne classification)
- Donner une sortie PDF qui résume les résultats des tests d'algorithmes

2.3 Contraintes de l'application

En ce qui concerne notre application web, différentes contraintes inhérentes à sa nature sont prises en compte à savoir:

- L'application doit être compatible avec n'importe quel jeu de données (format TSV ou CSV)
- L'application doit donc être réactive
- L'utilisateur doit pouvoir choisir les variables explicatives et les variables à expliquer
- L'utilisateur choisit quels algorithmes tester sur son jeu de données
- L'application doit être facile à maintenir.

2.4 Présentation générale

L'interface utilisateur a été conçue pour être intuitive et fonctionnelle. Répondant aux contraintes, elle permet à l'utilisateur de charger un jeu de données (de type TSV ou CSV), de choisir les variables explicatives et celles à expliquer, ainsi que les algorithmes et les paramètres de son choix. Outre cela, l'utilisateur a aussi la possibilité de faire des statistiques descriptives pour lui permettre d'explorer son jeu de données.

Ci-dessous une capture d'écran de l'application avec un exemple à l'appui.

a. Préparation de données

Notre application commence par la préparation des données. Dans cette fenêtre, l'utilisateur pourra:

- télécharger des fichiers de données (.csv / .tsv à ce stade de développement)
- choisir la présence de "header", le séparateur et la citation
- et enfin voir un aperçu rapide des données après avoir cliqué sur le bouton "Upload".

b. Statistiques Descriptives

Après avoir téléchargé le fichier de données dans l'application, dans la deuxième fenêtre, nous pourrions voir une statistique descriptive automatisée.

Sur le premier graphique, nous pouvons voir la description des informations de base dans la classe de données d'entrée:

- *lignes*: nombre de lignes
- *colonnes*: nombre de colonnes
- *discrete_columns*: nombre de colonnes discrètes
- *colonnes_continues*: nombre de colonnes continues
- *all_missing_columns*: nombre de colonnes avec tout ce qui manque
- *total_missing_values*: nombre d'observations manquantes
- *complete_rows*: nombre de lignes sans valeurs manquantes
- *total_observations*: nombre total d'observations
- *memory_usage*: allocation de mémoire estimée en octets

Dans le plot du milieu, nous observons la fréquence des valeurs manquantes pour chaque entité. Il y a 4 groupes pour décrire la situation des valeurs manquantes.

- *Good* : quand la proportion des valeurs manquantes ≤ 0.05
- *OK* : $0.05 < \text{la proportion des valeurs manquantes} \leq 0.4$
- *Bad* : $0.4 < \text{la proportion des valeurs manquantes} \leq 0.8$
- *Remove* : $0.8 < \text{la proportion des valeurs manquantes} \leq 1$

Sur le graphique de droite, nous pouvons voir nos variables une par une. Si la variable se compose de données numériques, un boxplot sera généré, et dans le cas de données non numériques, nous verrons un barplot.

c. ML Parameters

A cette étape, après avoir téléchargé le fichier de données et effectué quelques statistiques descriptives, l'utilisateur a la possibilité de choisir des paramètres pour les algorithmes souhaités.

Dans un premier temps, l'utilisateur fait le choix entre la régression et la classification.

Rappel

La classification est la tâche de prédire une étiquette de classe discrète.
La régression consiste à prédire une quantité continue.

Ensuite, pour le choix des hyperparamètres, il peut:

- soit saisir les paramètres définis par les siens

- soit laisser le package *caret* de R régler les hyperparamètres */parameter tuning/*

****Réglage des hyperparamètres avec caret****

Avec l'application, un fichier .csv est fourni avec les valeurs par défaut pour le réglage des hyperparamètres qui est visible ci-dessous.

Comme on peut le voir:

- pour certains des algorithmes, une mise à centrée réduit sera proposée comme *preProcessing*.
- la métrique de formation a été choisie precision */Accuracy/*
- la qualité des paramètres sera contrôlée par une validation croisée (10 fois)
- pour les algorithmes *glmnet /Ridge et Lasso/* α et les valeurs possibles de λ sont définies en complément.

****Réglage des hyperparamètres à la main****

Lorsque l'utilisateur décide d'entrer ses propres paramètres, il doit choisir les algorithmes que l'on souhaite utiliser dans la partie intitulée "Algorithmes" et entrer les paramètres. Après exécution, dans une nouvelle fenêtre on peut voir le choix final des paramètres qui seront utilisés pour la partie principale de modélisation.

3 Algorithmes de Machine Learning

Afin de répondre au mieux à notre problématique nous avons fait le choix d'utiliser différents algorithmes pour analyser nos données.

kNN /k-ppv/

L'algorithme des kNN consiste à trouver les k observations x_i les plus proches de l'observation x' à classer. Ensuite, il faut définir $f(x')$ en fonction des réponses y_i des kNN. Pour la régression c'est la valeur moyenne et pour la classification le vote majoritaire. Les questions ouvertes pour cet algorithme est le choix du **critère de distance** et de la **valeur de k**.

- kNN pour la régression:

$$\tilde{f}(x) = \text{moyenne}(y_i | i \in N_k(x))$$

\Rightarrow approxime directement la fonction de régression $E[Y|X]$

Nature de l'approximation

1. Espérance \rightarrow moyenne empirique
2. Valeur ponctuelle \rightarrow voisinage (dans le conditionnement)

- kNN pour la classification:

$$\tilde{f}(x) = \text{majorite}(y_i | i \in N_k(x)) = \arg \max_{l=1, \dots, K} \tilde{P}_l = \frac{1}{k} \sum_{i \in N_k(x)} 1(y_i = l)$$

\Rightarrow approxime directement le classifieur de Bayes:

$$\arg \max_{l=1, \dots, K} P(Y = C_l | X = x)$$

Nature de l'approximation

1. Probabilité \rightarrow proportion empirique
2. Valeur ponctuelle \rightarrow voisinage (dans le conditionnement)

CART - Classification et Regression Trees

De base, les algorithmes d'arbres de décision ne sont rien d'autre que des instructions if-else qui peuvent être utilisées pour prédire un résultat basé sur des données. La méthodologie CART ou Arbres de classification et de régression fait référence à ces deux types d'arbres de décision.

- arbres de classification

Un arbre de classification est un algorithme dans lequel la variable cible est fixe ou catégorielle. L'algorithme est ensuite utilisé pour identifier la «classe» dans laquelle une variable cible se situerait le plus probablement. Un arbre de classification divise l'ensemble de données en fonction de l'homogénéité des données. Disons, par exemple, qu'il y a deux variables; revenu et Âge; qui déterminent si un consommateur achètera ou non un type particulier de téléphone.

Si les données de formation montrent que 95% des personnes de plus de 30 ans ont acheté le téléphone, les données y sont divisées et l'Âge devient un noeud supérieur dans l'arbre. Cette division rend les données «pures à 95%». Des mesures d'impuretés comme l'entropie ou l'indice de Gini sont utilisées pour quantifier l'homogénéité des données en ce qui concerne les arbres de classification.

- arbres de régression

Un arbre de régression fait référence à un algorithme dans lequel se trouve la variable cible et l'algorithme utilisé pour prédire sa valeur. Dans un arbre de régression, un modèle de régression est ajusté à la variable cible en utilisant chacune des variables indépendantes. Après cela, les données sont divisées en plusieurs points pour chaque variable indépendante.

A chacun de ces points, l'erreur entre les valeurs prédites et les valeurs réelles est mise au carré pour obtenir une «somme des erreurs au carré» (SEC). Cette somme est comparée entre les variables. Celle ou le point qui a le SEC le plus bas est choisi comme point de partage. Ce processus se poursuit récursivement.

rfRanger ???

Random Forest

À la base, la méthode des forêts aléatoires est basée sur le bagging /aggrégation par bootstrap/, il est très performant sur de nombreux problèmes et facile à paramétrer. Principe d'algorithmes et de découper l'espace d'entrée en région, estimer et prédire une valeur par région. Pour la régression on choisit une valeur moyenne dans la région, et pour la classification on choisit la classe majoritaire. Les avantages de cette méthode sont l'interprétabilité du modèle et le mécanisme de prédiction proche du processus humain.

XGBoost

La bibliothèque XGBoost implémente l'algorithme d'arbres de décision de renforcement du gradient. Le boosting est une technique d'ensemble où de nouveaux modèles sont ajoutés pour corriger les erreurs commises par les modèles existants. Les modèles sont ajoutés séquentiellement jusqu'à ce qu'aucune autre amélioration ne puisse être apportée.

L'amplification du gradient est une approche où de nouveaux modèles sont créés et qui prédisent les résidus ou les erreurs des modèles précédents, puis additionnés pour faire la prédiction finale. Il est appelé boosting de gradient car il utilise un algorithme de descente de gradient pour minimiser la perte lors de l'ajout de nouveaux modèles.

Cette approche prend en charge à la fois les problèmes de modélisation prédictive de régression et de classification.

glmNetL2 - Ridge Regression

Régression linéaire avec perte quadratique et pénalité L2 :

$$(w^*, b^*) = \underset{w \in \mathbb{R}^1, b \in \mathbb{R}}{\operatorname{argmin}} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|w\|_2^2$$

$$= \underset{w \in \mathbb{R}^1, b \in \mathbb{R}}{\operatorname{argmin}} \sum_{i=1}^n (y_i - b - \sum_{j=1}^p w_j x_{ij} + \lambda \sum_{j=1}^p w_j^2)$$

la régression ridge impose donc une contrainte sur les coefficients (w). Le terme de pénalité (λ) régularise les coefficients de telle sorte que si les coefficients prennent de grandes valeurs, la fonction d'optimisation est pénalisée. Ainsi, la régression ridge réduit les coefficients et contribue à réduire la complexité du modèle et la multi-colinéarité.

Pour le résoudre il faut centrer les variables $\rightarrow \tilde{x}_{ij} = x_{ij} - \bar{x}$

1. on estime alors l'intercept b par $b^* = \hat{y} = \frac{1}{n} \sum_{i=1}^n y_i$
2. on obtient w avec le meme probleme sans intercept
3. solution : $w^* = (X^T X + \lambda I)^{-1} X^T y$ ou $X[i, j] = \tilde{x}_{ij}$

glmNetL1 - Lasso Regression

Régression linéaire avec perte quadratique et pénalité L1 :

$$\begin{aligned} (w^*, b^*) &= \underset{w \in \mathbb{R}^1, b \in \mathbb{R}}{\operatorname{argmin}} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|w\|_1 \\ &= \underset{w \in \mathbb{R}^1, b \in \mathbb{R}}{\operatorname{argmin}} \sum_{i=1}^n (y_i - b - \sum_{j=1}^p w_j x_{ij} + \lambda \sum_{j=1}^p |w_j|) \end{aligned}$$

Si on regarde l'équation en détail, on peut voir que la seule différence est qu'au lieu de prendre le carré des coefficients, les grandeurs sont prises en compte. Ce type de régularisation (L1) peut conduire à des coefficients nuls, c'est-à-dire que certaines caractéristiques sont complètement négligées pour l'évaluation des résultats. Ainsi, la régression Lasso aide non seulement à réduire le sur-ajustement, mais elle peut également nous aider à sélectionner les variables.

Du coup, par rapport à la pénalité ridge :

- même effet de régularisation : shrinkage des coefficients
- mais conduit à des coefficients exactement = 0

\Rightarrow solution parcimonieuse (sparse) : sélection de variables

AdaBoost Classification Trees ?????

AdaBoost M1

AdaBoost (boosting adaptatif) est un algorithme d'apprentissage assembliste qui peut être utilisé pour la classification ou la régression. Bien qu'AdaBoost résiste mieux au surajustement qu'un grand nombre d'algorithmes de machine learning, il est parfois sensible aux données bruitées et aux valeurs aberrantes.

AdaBoost est désigné comme adaptatif car il utilise de nombreuses itérations pour générer un apprenant composite fort. AdaBoost crée l'apprenant fort (un classifieur bien corrélé au classifieur correct) en ajoutant de manière itérative des apprenants faibles (un classifieur légèrement corrélé au classifieur correct). Lors de chaque session d'apprentissage, un nouvel apprenant faible est ajouté à l'ensemble et un vecteur pondération est ajusté afin de mettre l'accent sur les exemples ayant été classés de manière incorrecte lors des sessions précédentes. Par conséquent, le classifieur obtenu est doté d'une meilleure précision que classifieurs des apprenants faibles.

AdaBoost.M1 représente les algorithmes originaux de classification binaire.

Optimisation du modèle - Validation Croisée

L'optimisation du modèle signifie l'optimisation hyperparamétrique, c'est-à-dire il faut trouver le bon niveau de complexité du modèle, qui est en général difficile de choisir à priori. Pour l'optimisation on utilise la méthode de la validation croisée.

Au préalable, on définit un jeu de test pour évaluer les performances du modèle. Ensuite, on découpe le jeu d'apprentissage en K parties - les folds.

pour $k = 1$ à K , il faut :

- mettre de côté la $k - ième$ fold
- apprendre le modèle sur les $K - 1$ folds restantes
- appliquer le modèle sur les données de la $k - ième$ fold

A la fin, on évalue les performances du modèle.