

# Rapport Projet industriel : Prédiction de durées de séjours aux Hospices Civils de Lyon (HCL)

Grégoire MASSOT

25 Janvier 2016



## Table des matières

<b>1</b>	<b>Présentation du projet</b>	<b>2</b>
1.1	Contexte du projet . . . . .	2
1.2	Outils utilisés pour le projet . . . . .	2
1.3	Présentation de la base de données . . . . .	2
1.4	Chargement de la BDD . . . . .	4
1.5	Mesure de la performance des modèles . . . . .	4
<b>2</b>	<b>Modélisation du problème</b>	<b>5</b>
2.1	Présentation de la méthode . . . . .	5
2.2	Implémentation de la méthode . . . . .	5
2.2.1	Réglage des paramètres . . . . .	5
2.2.2	Construction du modèle . . . . .	5
<b>3</b>	<b>Analyse des résultats</b>	<b>6</b>
3.1	Importance des différentes variables . . . . .	6
3.1.1	1er Modèle . . . . .	6
3.1.2	2nd Modèle . . . . .	7
3.1.3	3ème Modèle . . . . .	7
3.2	Performance des modèles . . . . .	7
3.3	Incertitudes sur les résultats . . . . .	8

---

# 1 Présentation du projet

## 1.1 Contexte du projet

Ce projet industriel est orienté Data Science. Il s'agit de prédire le plus précisément possible des durées de séjour en milieu hospitalier à partir des données relatives aux patients à leur entrée à l'hôpital et des statistiques sur les patients précédents. Il s'agit donc d'un exercice d'apprentissage supervisé.

Ce projet s'inscrit dans une volonté de réduction des coûts des Hopitaux Civils de Lyon (HCL). Une maximisation de l'occupation des lits hospitaliers permet d'accueillir plus de patients avec les moyens disponibles.

## 1.2 Outils utilisés pour le projet

On utilise exclusivement R et son interface Rstudio pour établir les prédictions. On utilisera les bibliothèques utiles pour le problème, en particulier **xgboost**

On utilisera LaTeX pour écrire les rapports (comme celui-ci). On utilise Antidote pour la correction orthographique.

## 1.3 Présentation de la base de données

On travaille sur la base de données hospitalière `base_ano.txt` ayant le squelette suivant :

Duree du séjour	Âge du patient	Diagnostic principal	Service hospitalier	...
<b>duree1</b>	age1	dp1	service1	...
<b>duree2</b>	age2	dp2	service2	...
<b>duree3</b>	age3	dp3	service3	...
...	...	...	...	...

La table contient les prédictors suivants :

Prédicteur	Nature	Description	Modèle 1	Modèle 2	Modèle 3
<b>moisent</b>	Discret	Numéro du mois lors de l'entrée du patient	Non	Oui	Oui
<b>moissor</b>	Discret	Numéro du mois lors de la sortie du patient	Non	Non	Non
<b>age</b>	Continu	Age du patient, en années	Oui	Oui	Oui
<b>codepos</b>	Discret	Code postal du domicile du patient	Non	Oui	Oui
<b>MI-DEPR2</b>	Discrets	Comorbidités diagnostiquées	Non	Non	Oui
<b>Elix</b>	Discret	Nombre de comorbidités diagnostiquées	Non	Non	Oui
<b>centre</b>	Discret	Numéro du centre des HCL ou le patient a été pris en charge	Non	Oui	Oui
<b>service</b>	Discret	Numéro de l'unité médicale des HCL ou le patient a été pris en charge	Oui	Oui	Oui
<b>entree</b>	Discret	Mode d'entrée du patient	Non	Oui	Oui
<b>sortie</b>	Discret	Mode de sortie du patient	Non	Non	Non
<b>diag</b>	Discret	Diagnostic du patient selon la classification CIM-10	Oui	Oui	Oui
<b>ID</b>	Discret	ID du séjour	Non	Non	Oui

Dans le premier modèle, on garde seulement les prédicteurs de l'étude précédente. Dans le second modèle, on ajoute tous les prédicteurs qui sont connus de façon certaine lors de l'arrivée du patient.

Dans le dernier troisième modèle, on ajoute :

- Les comorbidités : elles ne sont pas connues de façon certaine lors de l'arrivée du patient
- l'ID : Il permet de situer temporellement le séjour par rapport aux autres dans la table. Ce n'est pas un prédicteur utilisable en pratique. Il permet cependant de mesurer l'importance de pouvoir situer les séjours dans le temps.

Le mois de sortie et le mode de sortie sont évidemment déterminés le jour de sortie du patient. On élimine donc ces données de la table pour tous les modèles.

## 1.4 Chargement de la BDD

On garde seulement les séjours dont la durée est supérieure ou égale à 1 jour. On découpe le jeu de données `donnees_hcl` en un `TrainData` sur lequel on construira le modèle et un `TestData` sur lequel on comparera les durées de séjours prédites et les durées de séjour réelles.

```
1 # Chargement de la BDD
2 donnees_hcl <- read.csv2(file = "base_ano.txt", header = TRUE, sep="
  \t")
3 # Sélection des variables (code à modifier en fonction du modèle)
4 donnees_hcl <- donnees_hcl[donnees_hcl$duree >= 1,]
5 donnees_hcl <- donnees_hcl[, -seq(8, 57)]
6 donnees_hcl <- subset(donnees_hcl, select = -c(moissor, sortie))
7
8 # Transformation du facteur "age" en variable continue
9 donnees_hcl$age <- as.character(donnees_hcl$age)
10 donnees_hcl$age <- as.numeric(donnees_hcl$age)
11
12 # Création des jeux de Train et de Test
13 TrainData <- donnees_hcl[donnees_hcl$Selected == 1,]
14 TestData <- donnees_hcl[donnees_hcl$Selected == 0,]
15
16 predcteurs_Train <- subset(TrainData, select = -c(duree))
17 duree_Train <- subset(TrainData, select = c(duree))
18
19 predcteurs_Test <- subset(TestData, select = -c(duree))
20 duree_Test <- subset(TestData, select = c(duree))
```

## 1.5 Mesure de la performance des modèles

On utilisera le critère RMSE (*Root Mean Square Error*) pour pouvoir mesurer la performance des modèles proposés. C'est le critère qui a été utilisé dans l'étude précédente. Cela permettra de mesurer la performance de cette étude.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - y_i^{pred})^2}$$

```
1 RMSE <- function(y, pred)
2 {
3   sqrt(mean((y - pred)^2))
4 }
```

On utilise aussi le critère MAE (*Mean Absolute Error*), plus parlant que le RMSE. C'est l'erreur moyenne sur les prévisions des durées de séjours (en jours)

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - y_i^{pred}|$$

```
1 MAE <- function(y, pred)
2 {
3   mean(abs(y - pred))
4 }
```

---

## 2 Modélisation du problème

### 2.1 Présentation de la méthode

Le Gradient Boosting est une méthode de machine learning parmi les plus fréquemment utilisées pour les problèmes d'apprentissage supervisés. La bibliothèque xgboost est souvent utilisée dans les solutions gagnantes des challenges kaggle.

### 2.2 Implémentation de la méthode

#### 2.2.1 Réglage des paramètres

Plusieurs paramètres sont à ajuster pour la méthode du gradient boosting. Les plus importants sont les suivants :

- `nrounds` : c'est le nombre d'arbres à construire. Plus il est grand, plus la prédiction est exacte et plus le temps de simulation est long
- `maxdepth` : c'est la profondeur maximale des arbres de décision.
- `eta` : permet de régler le compromis biais/variance lors de l'apprentissage.

On utilise la fonction `train` du package `caret` pour pouvoir choisir au mieux ces paramètres.

```
1 # Tuning des paramètres de xgboost
2 library(caret)
3 library(foreach)
4 library(doParallel)
5 cl <- makeCluster(8)
6
7 tuneGrid <- expand.grid(max_depth = c(2,4,6,8,10),
8                        nrounds = 100,
9                        eta = c(0.01,0.1,0.5,1))
10
11 fitControl <- trainControl(method = "repeatedcv",
12                           number = 3,
13                           repeats = 1)
14
15 xvalXGB <- train(TrainData$duree ~ .,
16                data = TrainData,
17                method = "xgbTree",
18                tuneGrid = tuneGrid,
19                trControl = fitControl)
```

#### 2.2.2 Construction du modèle

Une fois ce réglage effectué, on peut lancer le gradient boosting sur le `TrainData` avec un nombre d'arbres le plus grand possible pour avoir le modèle le plus précis construit dans un temps raisonnable.

```
1 # Construction du modèle
2 library(xgboost)
3 param <- list("objective" = "reg:linear",
4              "eta"=0.5,
5              "max.depth"=4,
6              "nthread" = 8)
7
8 modelXgboost <- xgboost(data = as.matrix(predicteurs_Train)
9                        , label = as.matrix(duree_Train))
```

```

10         , params=param
11         , nrounds = 100)
12 predictionXGBoost <- predict(modelXgboost ,
13                               as.matrix(predicteurs_Test))
14
15 RMSE(duree_Test, predictionXGBoost)

```

### 3 Analyse des résultats

#### 3.1 Importance des différentes variables

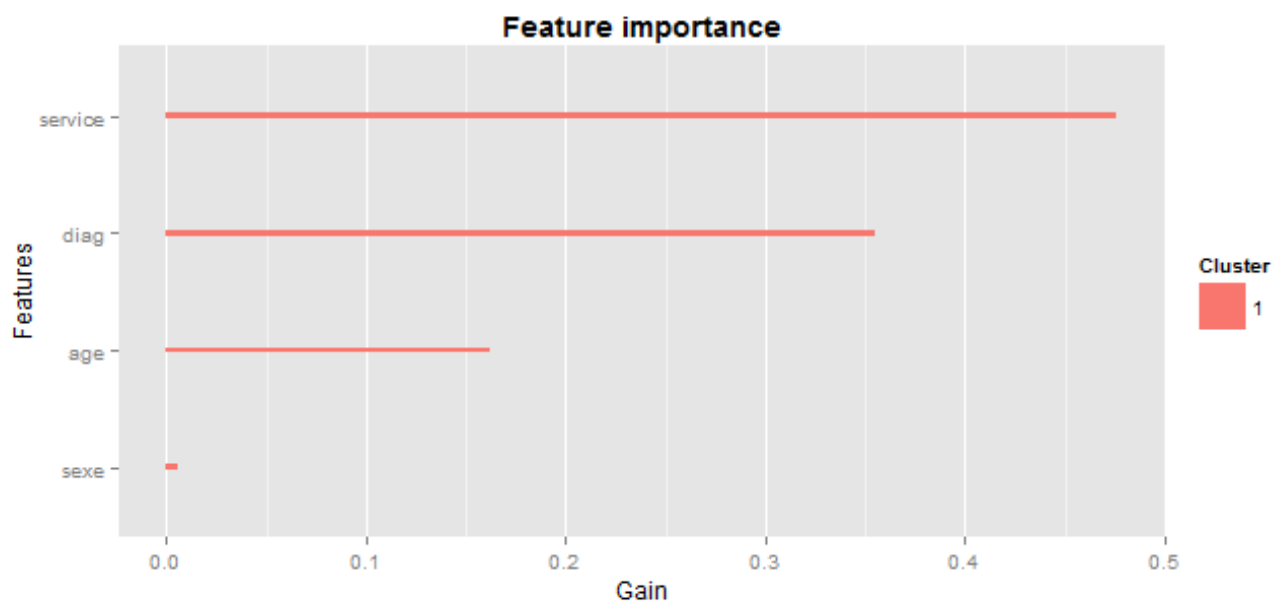
On utilise l'outil `xgb.importance` fourni avec `xgboost` pour déterminer les prédictors retenus par l'algorithme de machine learning.

```

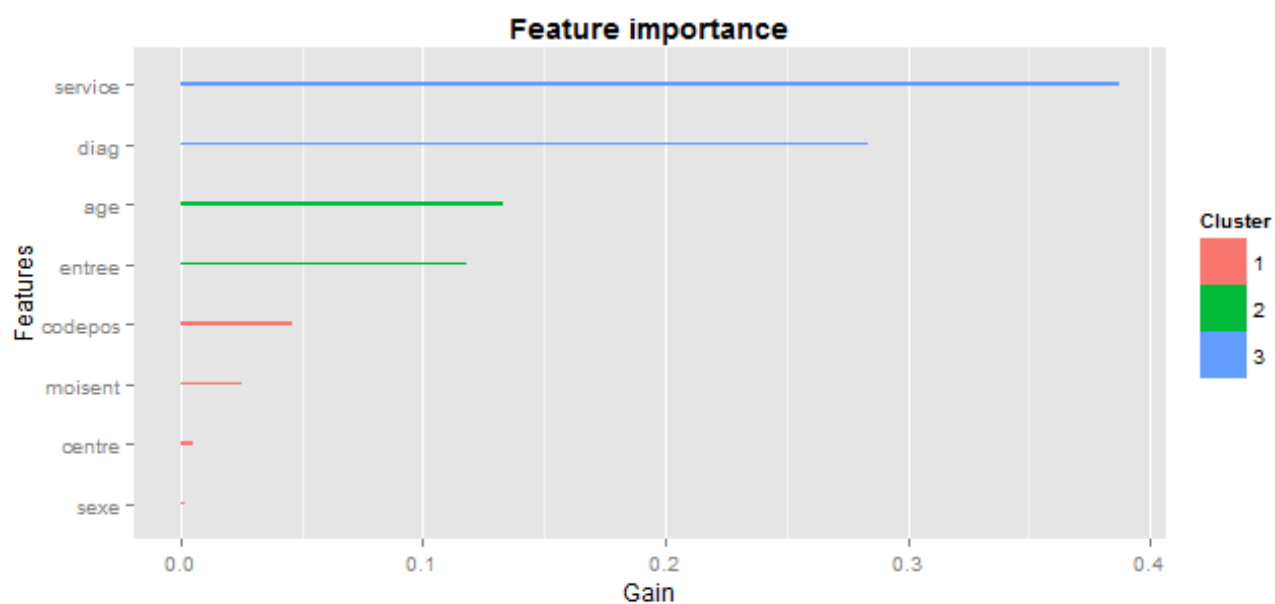
1 names <- dimnames(predicteurs_Train)[[2]]
2 importance_matrix <- xgb.importance(names, model = modelXgboost)
3 xgb.plot.importance(importance_matrix)
4 xgb.plot.tree(feature_names = names, model = modelXgboost, n_first_
  tree = 2)

```

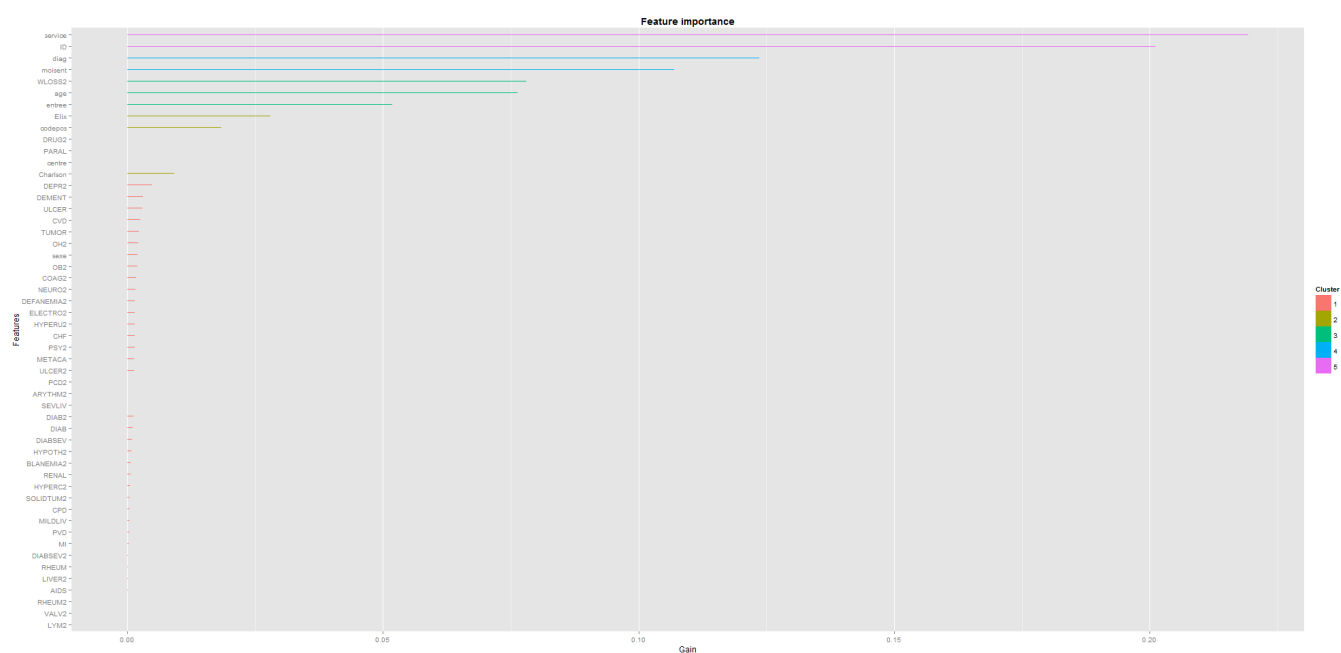
##### 3.1.1 1er Modèle



### 3.1.2 2nd Modèle



### 3.1.3 3ème Modèle



## 3.2 Performance des modèles

	RMSE	MAE
Modèle 1	6.88	3.14 jours
Modèle 2	6.84	3.15 jours
Modèle 3	5.68	2.75 jours

On voit que l'ajout de l'identifiant (qui permet de situer le séjour temporellement par rapport aux autres) ainsi que certaines comorbidités permettent de gagner environ 10% de MAE

### 3.3 Incertitudes sur les résultats

Pour améliorer la précision de la prédiction, on peut demander à xgboost de donner prédire non plus une durée brute, mais la densité de probabilité de la variable aléatoire relative à la durée de séjour.

Pour cela, on ajoute le code suivant lors du chargement des données

```
1 donnees_hcl <- donnees_hcl[donnees_hcl$duree >= 1 & donnees_hcl$duree <= 30 ,]
2 duree_Test_RMSE <- subset(donnees_hcl[donnees_hcl$Selected == 0,],
  select = c(duree))
3 donnees_hcl$duree <- as.factor(donnees_hcl$duree)
```

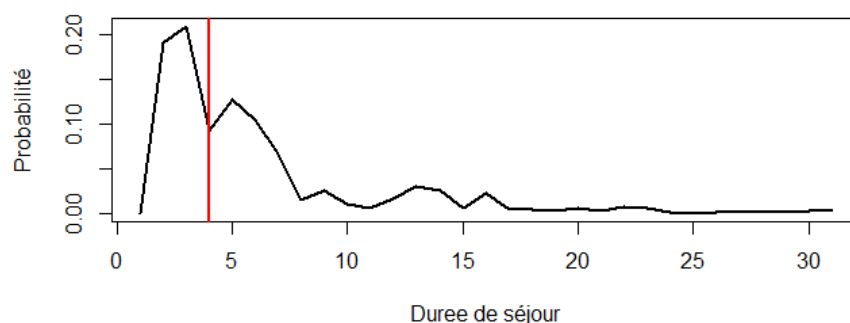
On remplace aussi les paramètres de xgboost.

```
1 param <- list("objective" = "multi:softprob",
2               "eta"=0.3,
3               "max.depth"=5,
4               "nthread" = 8,
5               "num_class" = 31)
6
7 modelXgboost <- xgboost(data = as.matrix(predicteurs_Train)
8                        , label = as.matrix(duree_Train)
9                        , params=param
10                       , nrounds = 100)
11 xgbtest <- xgb.DMatrix(data.matrix(predicteurs_Test), missing=NA)
12
13 xgbpred <- predict(modelXgboost, xgbtest)
14 probs <- t(matrix(xgbpred, nrow=31, ncol=length(xgbpred)/31))
```

On affiche ensuite les résultats, par exemple, les 100 premières entrées du jeu de test

```
1 for(i in 1:100)
2 {
3   plot(probs[i,], type='l')
4   abline(v=as.numeric(duree_Test$duree[i])+1, col="red")
5 }
```

On obtient des figures de ce type :





---

En noir la densité de probabilité de la variable aléatoire associée à la durée de séjour. En rouge la durée de séjour effectivement observée.

RMSE : 3.62 MAE : 2.24 jours

Les erreurs sont plus petites car on s'est limité aux séjours allant de moins de 30 jours afin de visualiser facilement les résultats sur un graphique.