



# Prévision de durées d'hospitalisation aux Hospices Civils de Lyon (HCL)

Grégoire MASSOT

21/01/2016





- **Entreprise** : Hospices civils de Lyon
- **Tuteur entreprise** : Antoine DUCLOS - Médecin à Lyon
- **Tuteur école** : Thierry GARAIX

# Déroulement du projet



4 venues à Lyon : 20 Octobre, 6 Novembre, 11 Décembre, 15 Janvier.

Les HCL doivent optimiser la gestion de leurs ressources pour accueillir le plus de patients avec les ressources disponibles.

Prévoir la durée d'hospitalisation des patients permet de maximiser l'occupation des lits

Les HCL souhaitent utiliser leur base de données sur les patients précédents pour prédire les durées d'hospitalisation des nouveaux patients.



**Objectif du PI** : déterminer la durée d'hospitalisation d'un patient à partir des données recueillies lors de son arrivée à l'hôpital et des statistiques sur les patients précédents.

Schéma de la Base de données anonymisée des HCL

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



Reprise du code R du challenge "Éolienne" de la majeure Data Science



Inscription au challenge Walmart sur kaggle.com

Grégoire MASSOT



Verified  
account

KAGGLER



Highest†  
**12765th**

Current†  
**13375th**  
/ 448,307



402.4 points  
Joined a month ago  
Til ranking method changed 13 May 2015 (?)

Profile

Results

Scripts

Forum

Account

Activity



**Walmart Recruiting: Trip Type Classification**

10 entries in team Grégoire MASSOT

Finished

**773rd/1047**



```
1 # Chargement de la BDD
2 donnees_hcl <- read.csv2(file = "base_ano.txt", header =
  TRUE, sep="\t")
3 # Sélection des séjours de durée supérieure ou égale à 1
  jour
4 donnees_hcl <- donnees_hcl[donnees_hcl$duree >= 1,]
5 # Sélection des prédicteurs effectivement disponibles
  lors de l'entrée du patient
6 donnees_hcl <- donnees_hcl[, -seq(8, 57)]
7 donnees_hcl <- subset(donnees_hcl, select = -c(moissor))
8 donnees_hcl <- subset(donnees_hcl, select = -c(sortie))
9 donnees_hcl <- subset(donnees_hcl, select = -c(ID))
10
11 # Transformation du facteur "age" en variable continue
12 donnees_hcl$age <- as.character(donnees_hcl$age)
13 donnees_hcl$age <- as.numeric(donnees_hcl$age)
```

```
1 # Création des jeux de Train et de Test
2 TrainData <- donnees_hcl[donnees_hcl$Selected == 1,]
3 TestData <- donnees_hcl[donnees_hcl$Selected == 0,]
4
5 predcteurs_Train <- subset(TrainData, select = -c(duree
6   ))
7
8 predcteurs_Test <- subset(TestData, select = -c(duree))
9 duree_Test <- subset(TestData, select = c(duree))
```



```
1 # Déclaration de la fonction RMSE
2
3 RMSE <- function(y, pred)
4 {
5   sqrt(mean((y - pred)^2))
6 }
```

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

```

1 # Tuning des paramètres de xgboost avec caret
2 library(caret)
3 library(foreach)
4 library(doParallel)
5 cl <- makeCluster(8)
6
7 tuneGrid <- expand.grid(max_depth = c(1,2,3,4,5,6,7,8,9)
8   ,
9   nrounds = 100,
10   eta = c
11     (0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,
12     )
13
14 fitControl <- trainControl(method = "repeatedcv",
15   number = 3,
16   repeats = 1)
17
18 xvalXGB <- train(TrainData$duree ~ .,
19   data = TrainData,
20   method = "xgbTree",
21   tuneGrid = tuneGrid,
22   trControl = fitControl)
23
24 xvalXGB

```

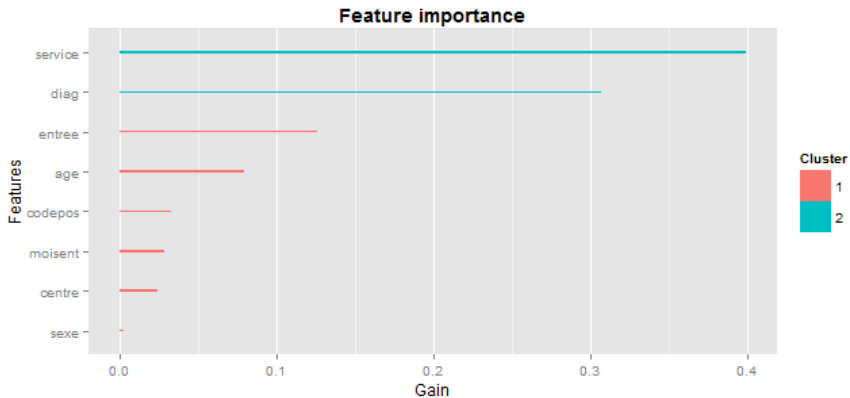
```
1 # Construction du modèle
2 library(xgboost)
3 param <- list("objective" = "reg:linear",
4               "eta"=0.5,
5               "max.depth"=5,
6               "nthread" = 8)
7
8 modelXgboost <- xgboost(data = as.matrix(predicteurs_
9                               Train)
10                          , label = as.matrix(duree_Train)
11                          , params=param
12                          , nrounds = 100)
13 predictionXGBoost <- predict(modelXgboost ,
14                               as.matrix(predicteurs_Test)
15                               )
16
17 RMSE(duree_Test, predictionXGBoost)
```

```
1 # feature importance
2
3 names <- dimnames(predicteurs_Train)[[2]]
4 importance_matrix <- xgb.importance(names, model =
    modelXgboost)
5 xgb.plot.importance(importance_matrix)
6 xgb.plot.tree(feature_names = names, model =
    modelXgboost, n_first_tree = 2)
```

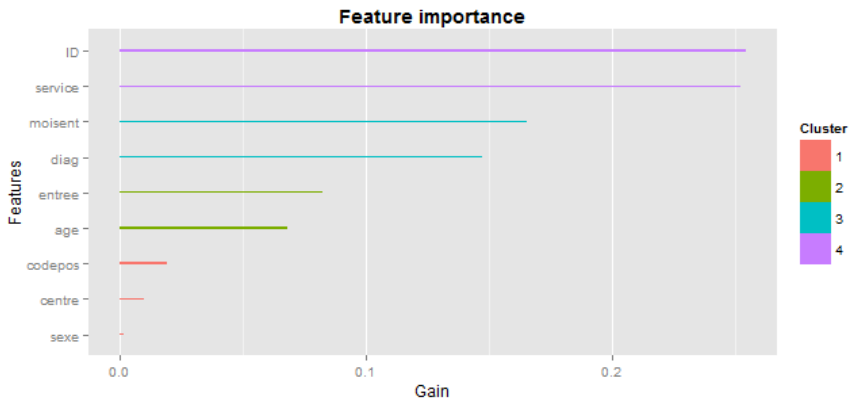


Méthode	RMSE
Régression linéaire (étude précédente menée par les HCL)	6.81
Machine learning - sans la variable ID	6.86
Machine learning - avec la variable ID	5.73

## Analyse de l'importance des prédicteurs - Sans ID



## Analyse de l'importance des prédicteurs - Avec ID





Fin