



## XGBoost, origines et applications

Damien DOUTEAUX

## Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

Sommaire

Aspects théoriques

Mise en œuvre

Applications

Conclusion

## Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

- Moteurs de recherche



- Alertes mails



- Réseaux sociaux et autres



Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

XGBoost : *EX*treme **G**radient **B**oosting

- ⊙ **Flexibilité** Régression, classification,...
- ⊙ **Portabilité** Windows, Linux, OS X
- ⊙ **Multi-langages** Python, R, JAVA, C++, Scala,...
- ⊙ **Distribué** Yarn, Spark, Flink, AWS, Azure,...
- ⊙ **Performance** Optimisé et expensif

Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

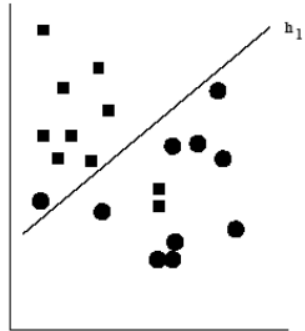
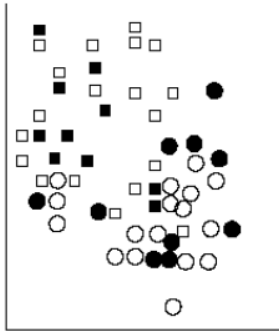
Conclusion

- ⊙ Une stratégie adaptative.
- ⊙ Convertir des règles peu performantes en (très) bonne prédiction.
- ⊙ Réduction variance et biais.
- ⊙ Convergence rapide.
- ⊙ Sensible au bruit.

# LE BOOSTING, PREMIER ALGORITHME

Sommaire  
Aspects  
théoriques  
Mise en  
œuvre  
Applications  
Conclusion

## PREMIER MODÈLE



Source : Cours Machine Learning, Haytham ELGHAZEL

# LE BOOSTING, PREMIER ALGORITHME

Sommaire

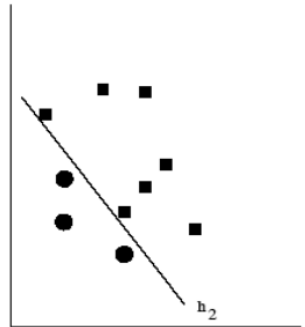
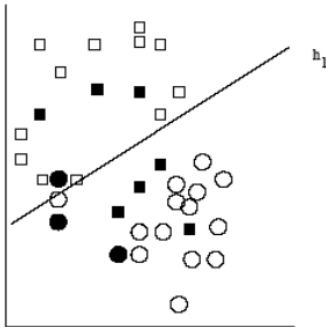
Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

## DEUXIÈME MODÈLE



Source : Cours Machine Learning, Haytham ELGHAZEL

# LE BOOSTING, PREMIER ALGORITHME

Sommaire

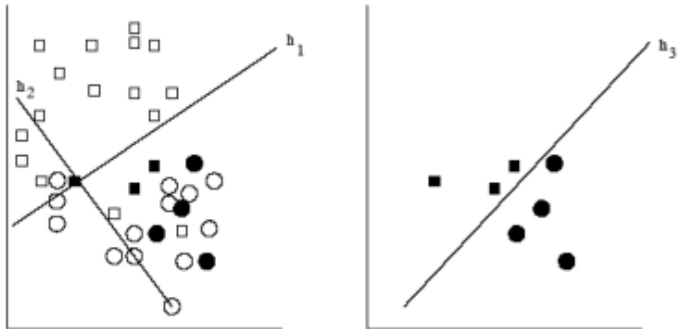
Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

## TROISIÈME MODÈLE



Source : Cours Machine Learning, Haytham ELGHAZEL



# LE BOOSTING, PREMIER ALGORITHME

Sommaire

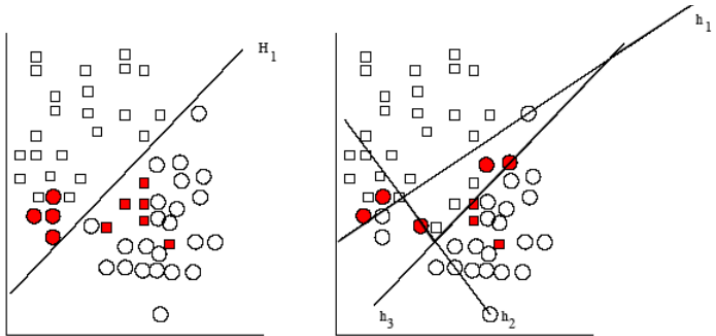
Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

## VOTE MAJORITAIRE



Source : Cours Machine Learning, Haytham ELGHAZEL

Sommaire

Aspects  
théoriques

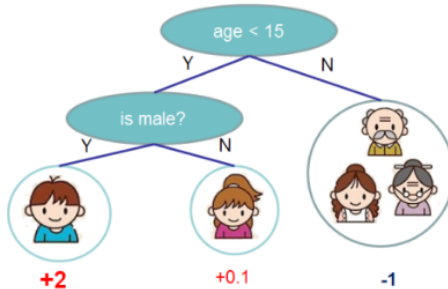
Mise en  
œuvre

Applications

Conclusion

## UN ARBRE SIMPLE (CART)

Does the person like computer games



Source : <https://xgboost.readthedocs.io/en/latest/model.html>

Sommaire

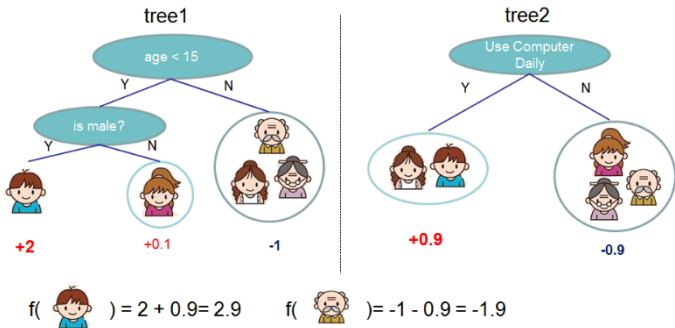
Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

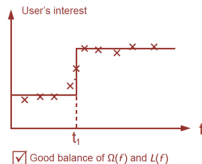
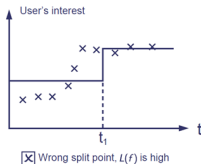
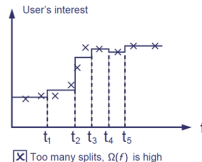
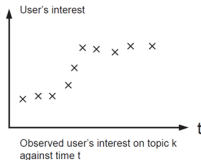
## PLUSIEURS ARBRES VALLENT MIEUX QU'UN



Source : <https://xgboost.readthedocs.io/en/latest/model.html>

## CHOIX DE L'ARBRE À AJOUTER

⊙ **Fonction objectif**  $\text{obj}(\Theta) = \mathcal{L}(\Theta) + \Omega(\Theta)$



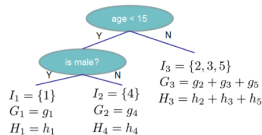
Source : <https://xgboost.readthedocs.io/en/latest/model.html>

## CHOIX DE L'ARBRE À AJOUTER

- ⊙ **Fonction de perte**  $\mathcal{L}(t) = \sum_{i=1}^n \left[ g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right]$
- ⊙ **Complexité**  $\Omega(t) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T \omega_j^2$

Instance index    gradient statistics

1		$g_1, h_1$
2		$g_2, h_2$
3		$g_3, h_3$
4		$g_4, h_4$
5		$g_5, h_5$



$$Obj = - \sum_j \frac{G_j^2}{H_j + \lambda} + 3\gamma$$

The smaller the score is, the better the structure is

Source : <https://xgboost.readthedocs.io/en/latest/model.html>

Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

## CHOIX DE L'ARBRE À AJOUTER

- ⦿ Énumérer toutes les structures d'arbres possibles.
- ⦿ Calculer l'objectif pour chaque structure.
- ⦿ Trouver l'optimal et optimiser les feuilles.

Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

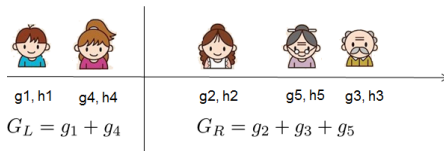
## CHOIX DE L'ARBRE À AJOUTER

- ⦿ Énumérer toutes les structures d'arbres possibles.
- ⦿ Calculer l'objectif pour chaque structure.
- ⦿ Trouver l'optimal et optimiser les feuilles.

**En pratique, construction des arbres au coup par coup.**

## CHOIX DE L'ARBRE À AJOUTER

⊙ Le gain  $\frac{1}{2} \left[ \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \gamma$



On ne s'arrête pas si  $\text{Gain} < 0$ .

Source : <https://xgboost.readthedocs.io/en/latest/model.html>



Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

- ⊙ Prise en compte de la régularisation.
- ⊙ Calcul en parallèle.
- ⊙ Support de Hadoop.
- ⊙ Possibilité d'adaptation des fonctions objectifs.
- ⊙ Prise en charge des valeurs manquantes.
- ⊙ Version améliorée de l'élagage
- ⊙ Cross-validation native

Sommaire

Aspects  
théoriques

Mise en  
œuvre

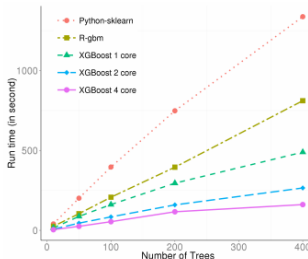
Applications

Conclusion

- ⊙ Entraîner les arbres.
- ⊙ Pour chaque variable :
  - ▷ Compter le nombre de fois où elle est sélectionnée.
  - ▷ Pondérer par la diminution d'erreur engendrée.
  - ▷ Moyenner sur les arbres.

La rapidité est le but initial de XGBoost :

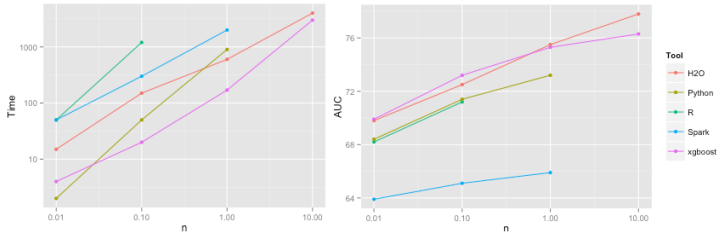
- ⊙ **Mémoire** Pas de mémoire dynamique.
- ⊙ **Cache** Éviter de surcharger la mémoire.
- ⊙ **Amélioration modèle** Voir précédemment.
- ⊙ **Conception** Parallélisation en arrière plan.
- ⊙ **Données externalisées** Si mémoire insuffisante.



Source : <http://www.jmlr.org/proceedings/papers/v42/chen14.pdf>

La rapidité est le but initial de XGBoost :

- ⊙ **Mémoire** Pas de mémoire dynamique.
- ⊙ **Cache** Éviter de surcharger la mémoire.
- ⊙ **Amélioration modèle** Voir précédemment.
- ⊙ **Conception** Parallélisation en arrière plan.
- ⊙ **Données externalisées** Si mémoire insuffisante.



Source : <http://datascience.la/benchmarking-random-forest-implementations/> [2015]

Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

## SUR LE BOOSTING

---

- ⊙ **1989** Boosting (R. Schapire)
- ⊙ **1996** AdaBoost (Y. Freund et R. Schapire)
- ⊙ **1999** GBM (L. Breiman puis J. Friedman)
- ⊙ **2014** XGBoost (T. Chen)

## POUR XGBOOST

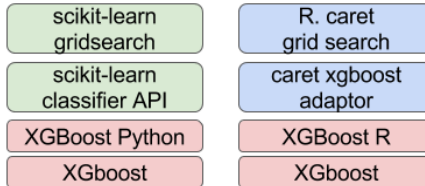
- ⊙ **Mars 2014** Première release
- ⊙ **Mai 2014** Python



Source : <http://homes.cs.washington.edu/~tqchen/2016/03/10/story-and-lessons-behind-the-evolution-of-xgboost.html>

## POUR XGBOOST

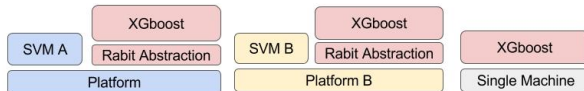
- ⊙ **Mars 2014** Première release
- ⊙ **Mai 2014** Python
- ⊙ **Septembre 2014** Parallélisation, R
- ⊙ **Mai 2015** YARN, gestion HDFS, SKLearn wrapper



Source : <http://homes.cs.washington.edu/~tqchen/2016/03/10/story-and-lessons-behind-the-evolution-of-xgboost.html>

## POUR XGBOOST

- ⊙ **Mars 2014** Première release
- ⊙ **Mai 2014** Python
- ⊙ **Septembre 2014** Parallélisation, R
- ⊙ **Mai 2015** YARN, gestion HDFS, SKLearn wrapper
- ⊙ **Janvier 2016** API JAVA, Flink, Spark, améliorations
- ⊙ **Juillet 2016** Totale compatibilité JVM (Spark,...)



Source : <http://homes.cs.washington.edu/~tqchen/2016/03/10/story-and-lessons-behind-the-evolution-of-xgboost.html>



Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

## PARAMÈTRES GÉNÉRIQUES

Pour définir par exemple quelle méthode Boosting sera utilisée.

## PARAMÈTRES LIÉS AU BOOSTING

Pour paramétrer le booster choisi.

## PARAMÈTRES LIÉS À L'APPRENTISSAGE

Dépend de la tâche d'apprentissage (classification,...).

Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

- ⊙ **Booster** Linéaire ou arbre.
- ⊙ **Silent** Affichage de messages.
- ⊙ **Nthread** Par défaut le maximum possible.

Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

*Pour celui sur les arbres. Douze paramètres utiles...*

- ⊙ **Eta**    Contrôle du niveau d'apprentissage.
- ⊙ **Min\_child\_weight**    Pour contrôler l'over/under-fitting
- ⊙ **Max\_depth**    Pour contrôler l'over-fitting.
- ⊙ **Lambda**    Pour de la régularisation.
- ⊙ **Gamma**    Valeur minimale de gain pour diviser.
- ⊙ ...

Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

- ⊙ **Objective** Fonction objectif à minimiser (linéaire, softmax, softprob,...).
- ⊙ **Eval\_metric** Métrique d'évaluation (erreur MSE, MAE, LogLoss, AUC,...).
- ⊙ **Seed** Pour l'aléatoire.

## Sommaire

## Aspects théoriques

## Mise en œuvre

## Applications

## Conclusion

```
1 val spark = SparkSession.builder().appName("SimpleXGBoost Application").
    config("spark.executor.memory", "2G").config("spark.executor.cores", "4")
    .config("spark.default.parallelism", "4").master("local[*]").
    getOrCreate()

2
3 // number of iterations
4 val numRound = 10
5 val numWorkers = 4
6 // training parameters
7 val paramMap = List(
8     "eta" -> 0.023f,
9     "max_depth" -> 10,
10    "min_child_weight" -> 3.0,
11    "subsample" -> 1.0,
12    "colsample_bytree" -> 0.82,
13    "colsample_bylevel" -> 0.9,
14    "base_score" -> 0.005,
15    "eval_metric" -> "auc",
16    "seed" -> 49,
17    "silent" -> 1,
18    "objective" -> "binary:logistic").toMap
19 println("Starting Xgboost ")
20 val xgBoostModelWithDF = XGBoost.trainWithDataFrame(trainingData, paramMap,
    round = numRound, nWorkers = numWorkers, useExternalMemory = true)
21
22 val predictions = xgBoostModelWithDF.setExternalMemory(true).transform(
    testData).select("label", "probabilities")
```

Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

```
1 loglossobj <- function(preds, dtrain) {  
2   # dtrain is the internal format of the training data  
3   # We extract the labels from the training data  
4   labels <- getinfo(dtrain, "label")  
5   # We compute the 1st and 2nd gradient, as grad and hess  
6   preds <- 1/(1 + exp(-preds))  
7   grad <- preds - labels  
8   hess <- preds * (1 - preds)  
9   # Return the result as a list  
10  return(list(grad = grad, hess = hess))  
11 }  
12  
13 model <- xgboost(data = train$data, label = train$label,  
14                 nrounds = 2, objective = loglossobj, eval_metric = "error")
```

# SÉLECTION DE VARIABLES

Sommaire

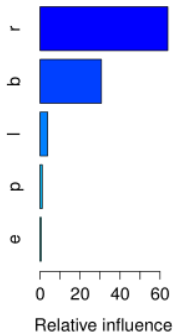
Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

```
1 bst <- xgboost(data = train$data, label = train$label, max.depth = 2,  
2               eta = 1, nthread = 2, nround = 2, objective = "binary:logistic")  
3 importance_matrix <- xgb.importance(agaricus.train$data@Dimnames[[2]], model  
   = bst)  
4 xgb.plot.importance(importance_matrix)
```



Source : <http://dmlc.ml/rstats/2016/03/10/xgboost.html> et cours d'apprentissage statistique (C. HELBERT)

Vendredi 3 mars 2017

18

Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

- ⊙ **1.** Fixer un niveau d'apprentissage élevé.
- ⊙ **2.** Trouver le nombre optimal d'arbres.
- ⊙ **3.** Gérer les paramètres des arbres.
- ⊙ **4.** Gérer les paramètres de régularisation.
- ⊙ **5.** Réduire le niveau d'apprentissage.
- ⊙ **6.** Utiliser l'AUC pour estimer les modèles.



Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

*En 2015 sur Kaggle, 17 solutions gagnantes sur 29  
utilisaient XGBoost.*

*En 2015 sur Kaggle, 17 solutions gagnantes sur 29 utilisaient XGBoost.*

- ⊙ **1<sup>er</sup>** Knowledge Discovery and Data Mining Cup 2016 (V. Sandulescu).
- ⊙ **1<sup>er</sup> et 3<sup>ème</sup>** CERN LHCb experiment Flavour of Physics competition 2015 (V. Mironov).
- ⊙ **1<sup>er</sup>** Caterpillar Tube Pricing competition (M. Filho).
- ⊙ **2<sup>ème</sup>** Airbnb New User Bookings (K. Kuroyanagi).
- ⊙ **2<sup>ème</sup>** Allstate Claims Severity (A. Noskov).
- ⊙ **10%** Higgs Boson Competition (T. Chen).
- ⊙ ...

# UN EXEMPLE CONCRET

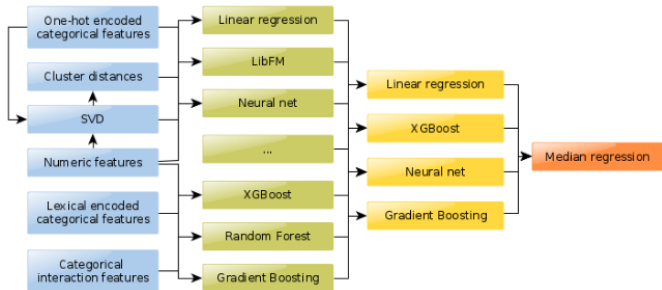
Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion



Source : <http://blog.kaggle.com/2017/02/27/allstate-claims-severity-competition-2nd-place-winners-interview-alezey-noskov/>

Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

*Des données difficiles à obtenir...*

- ⦿ ODPS Cloud Service (Alibaba)
- ⦿ Tencent (QQ)
- ⦿ AutoHome
- ⦿ AXA, Expedia, Amazon,...

Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

- ⦿ Une implémentation récente (3 ans).
- ⦿ Une forte portabilité et de bonnes performances.
- ⦿ Une utilisation industrielle qui se développe.
- ⦿ Savoir-faire nécessaire pour la configuration.

Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

- ⦿ Une implémentation récente (3 ans).
- ⦿ Une forte portabilité et de bonnes performances.
- ⦿ Une utilisation industrielle qui se développe.
- ⦿ Savoir-faire nécessaire pour la configuration.

**Une solution qui semble avoir de l'avenir!**

Sommaire

Aspects  
théoriques

Mise en  
œuvre

Applications

Conclusion

# Merci pour votre attention

# Et place aux questions!

 @DDouteaux

[https://ddouteaux.github.io/XGBoost\\_Veille\\_Douteaux/index.html](https://ddouteaux.github.io/XGBoost_Veille_Douteaux/index.html)