

# Pratique de l'apprentissage statistique

## 7. Boosting

V. Lefieux



École des Ponts

ParisTech

# Plan

Introduction

AdaBoost

Gradient boosting

Fonctions de pertes classiques

# Plan

Introduction

AdaBoost

Gradient boosting

Fonctions de pertes classiques

# Principes

- ▶ Le boosting est une méthode d'agrégation de prédicteurs récursive : le prédicteur obtenu à l'itération  $b$  dépend de celui obtenu à l'itération  $b - 1$ .
- ▶ Le boosting consiste à utiliser une **règle faible** qui apprend tout d'abord sur l'échantillon le plus simple, celui de départ, puis au fur et à mesure sur des échantillons rendus plus complexes par des pondérations adéquates, récursives.
- ▶ Contrairement au bagging, le boosting peut conduire à du sur-apprentissage : il faudra donc veiller à limiter le nombre d'itérations.
- ▶ C'est une méthode souvent utilisée en pratique au vu des bons résultats obtenus.
- ▶ Le boosting permet de traiter des problématiques de régression et de classification supervisée.

## Données considérées

- ▶ On dispose d'un échantillon de  $(X, Y)$  :

$$\mathcal{D}_n = (X_i, Y_i)_{i \in \{1, \dots, n\}} .$$

On note :

$$d_n = (x_i, y_i)_{i \in \{1, \dots, n\}} .$$

- ▶ On considère dans la suite que :

- ▶  $X \in \mathbb{R}^p$  :

*Toutes les covariables sont considérés quantitatives.*

*Mais il est également possible de considérer des covariables qualitatives.*

- ▶  $Y \in \{-1, 1\}$  dans le cas de la **classification supervisée** :

*On se place dans le cadre d'une classification supervisée binaire*

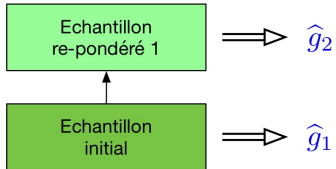
*Mais il est également possible de considérer des classifications supervisées avec  $K$  classes.*

- ▶  $Y \in \mathbb{R}$  dans le cas de la **régression**.

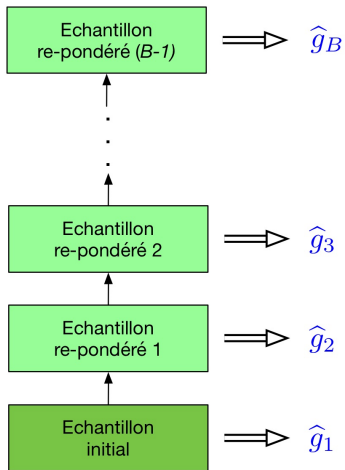
# Principe I



## Principe II

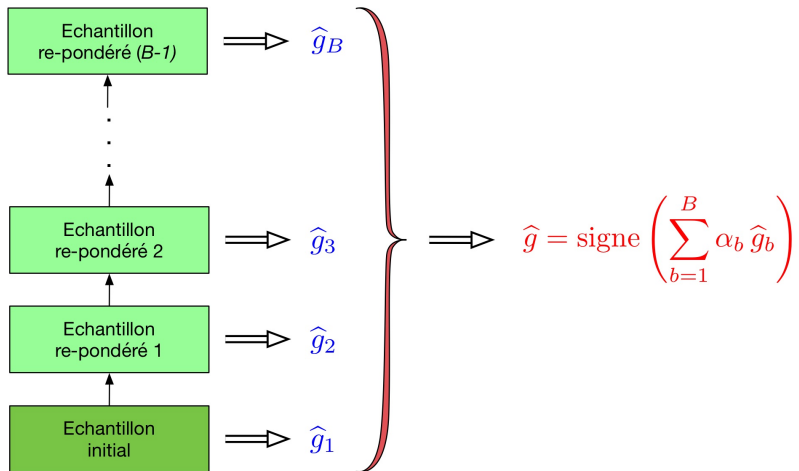


# Principe III





## Principe IV



# Enjeux

- ▶ Comment **repondérer les observations** à chaque étape ?
- ▶ Quels **poids accorder aux estimateurs** à chaque étape ?

# Plan

Introduction

AdaBoost

Gradient boosting

Fonctions de pertes classiques

# Règle faible

- ▶ On appelle une **règle faible** (*weak learner*) un prédicteur légèrement meilleur que le hasard.
- ▶ Par exemple, dans le cas de classification :

$$\mathbb{P}(g(X) \neq Y) = \frac{1}{2} - \gamma .$$

où  $\gamma > 0$ .

- ▶ On considère usuellement comme règle faible : 1-plus proche voisin, arbre à 2 feuilles, etc.

# Historiquement : AdaBoost (Freund et Schapire, 1997) I

## 1. Initialisation des poids des individus :

$$\forall i \in \{1, \dots, n\} : w_i^{(1)} = \frac{1}{n} .$$

## 2. Pour $b \in \{1, \dots, B\}$ :

2.1 Estimer  $\hat{g}_b$  avec les poids  $(w_1^{(b)}, \dots, w_n^{(b)})$  pour l'échantillon.

2.2 Calcul du taux d'erreur  $e_b$  de  $\hat{g}_b$  :

$$e_b = \frac{\sum_{i=1}^n w_i^{(b)} \mathbb{1}_{\hat{g}_b(x_i) \neq y_i}}{\sum_{i=1}^n w_i^{(b)}} .$$

2.3 Calcul de la pénalité :

$$\alpha_b = \ln \left( \frac{1 - e_b}{e_b} \right) .$$

2.4 Calcul des nouveaux poids :

$$\forall i \in \{1, \dots, n\} : w_i^{(b+1)} = w_i^{(b)} \exp \left( \alpha_b \mathbb{1}_{\hat{g}_b(x_i) \neq y_i} \right) .$$

## Historiquement : AdaBoost (Freund et Schapire, 1997) II

3. Le prédicteur obtenu au final est :

$$\hat{g} = \text{signe} \left( \sum_{b=1}^B \alpha_b \hat{g}_b \right) .$$

## Remarques

- ▶ L'étape 2.1 implique que la méthode retenue soit en mesure prendre en compte des poids. Dans le cas contraire, l'étape 2.1 d'estimation de  $\hat{g}_b$  s'effectue sur un échantillon de dimension  $n$  issu du tirage au sort d'observations de  $d_n$  avec remise selon les poids  $(w_1^{(b)}, \dots, w_n^{(b)})$ .
- ▶ Les poids sont modifiés de manière à accroître l'importance des observations mal classées et diminuer celle des observations bien classées.
- ▶ Le poids  $\alpha_b$  du prédicteur  $\hat{g}_b$  augmente avec sa performance :  $\alpha_b$  augmente lorsque  $e_b$  diminue.

## Quelques résultats théoriques

- On a :

$$R_n(\hat{g}) \leq \exp \left( -2 \sum_{b=1}^B \gamma_b^2 \right)$$

où  $\gamma_b = \frac{1}{2} - e_b$  est le **gain** de  $\hat{g}_b$  par rapport à une décision basée sur le hasard pur.

Si on est meilleur que le hasard alors le risque empirique tend exponentiellement vers 0 avec  $B$ .

- On peut montrer que plus  $B$  est grand, plus le biais est faible mais plus la variance est élevée.  
Il y a un risque de sur-apprentissage si  $B$  est trop important. Il faut donc contrôler  $B$  par validation croisée.
- On peut montrer qu'AdaBoost est un algorithme qui permet de minimiser le risque empirique « convexifié », et s'inscrit ainsi dans le cadre des méthodes de **gradient boosting**.



# Plan

Introduction

AdaBoost

Gradient boosting

Fonctions de pertes classiques

## Une idée générale : la convexification du risque I

- ▶ Si on considère la fonction de perte suivante :

$$\ell(y, y') = \mathbb{1}_{y \neq y'} ,$$

le risque du prédicteur  $g$  vaut :

$$R(g) = \mathbb{E}(\mathbb{1}_{g(X) \neq Y}) = \mathbb{P}(g(X) \neq Y) .$$

- ▶ On aimerait trouver l'estimateur qui minimise ce risque, malheureusement ce dernier n'est pas calculable en pratique.
- ▶ On considère alors le risque empirique sur  $\mathcal{D}_n$  :

$$R_n(g) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{g(X_i) \neq Y_i} .$$

- ▶ Malheureusement, ce risque empirique n'est pas convexe et donc difficile à optimiser.

## Une idée générale : la convexification du risque II

- ▶ On souhaite donc trouver sur une fonction de perte telle que le risque empirique soit convexe, et donc facile à optimiser.
- ▶ Il nous suffit pour cela de considérer une fonction  $\ell : (u, y) \mapsto \ell(u, y)$  convexe en  $u$ .
- ▶ Dans le cas de la classification supervisée, on peut par exemple choisir :

$$\ell : \mathbb{R} \times \{-1, 1\} \rightarrow \mathbb{R}^+$$

$$(u, y) \mapsto \exp(-y u)$$

qui est bien convexe en  $u$ .

# Optimisation

- ▶ On considère donc une fonction de perte telle que le risque empirique soit convexe, et donc facile à optimiser.
- ▶ On recherche la solution, de manière récursive, au problème :

$$\min_g \frac{1}{n} \sum_{i=1}^n \ell(g(X_i), Y_i) .$$

- ▶ On peut utiliser l'algorithme de Newton-Raphson ou plutôt l'algorithme de descente de gradient fonctionnel.

# Algorithme de Newton-Raphson

- Soit :

$$J(g) = \sum_{i=1}^n \ell(g(x_i), y_i) .$$

- Avec la notation :

$$\widehat{g}_b = (\widehat{g}_b(x_1), \dots, \widehat{g}_b(x_n))^{\top} ,$$

la formule de récurrence de l'algorithme de Newton-Raphson est :

$$\widehat{g}_b = \widehat{g}_{b-1} - \lambda \nabla J(\widehat{g}_{b-1})$$

où  $\lambda$  est un paramètre de régularisation.

- La fonction  $\widehat{g}_b$  n'est calculées qu'aux points  $(x_1, \dots, x_n)$ , c'est notamment pour cela qu'on préfère utiliser l'algorithme de descente de gradient fonctionnel.

# Algorithme de descente de gradient fonctionnel : cas de la classification supervisée

## 1. Initialisation :

$$\widehat{g}_0 = \arg \min_{c \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \ell(c, y_i) .$$

## 2. Pour $b \in \{1, \dots, B\}$ :

2.1 Pour  $i \in \{1, \dots, n\}$  : calculer les opposés du gradient :

$$U_i = - \frac{\partial}{\partial g(x_i)} \ell(g(x_i), y_i) \Big|_{g(x_i) = \widehat{g}_{b-1}(x_i)} .$$

2.2 Estimer la règle faible sur l'échantillon  $((x_1, U_1), \dots, (x_n, U_n))$ .  
On obtient ainsi  $\widehat{h}_b$ .

2.3 Mettre à jour :

$$\widehat{g}_b(x) = \widehat{g}_{b-1}(x) + \lambda \widehat{h}_b(x)$$

où  $\lambda$  est un paramètre de régularisation.

## 3. Le prédicteur obtenu au final est :

$$\widehat{g} = \text{signe}(\widehat{g}_B) .$$

# Algorithme de descente de gradient fonctionnel : cas de la régression

## 1. Initialisation :

$$\hat{m}_0 = \arg \min_{c \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \ell(c, y_i) .$$

## 2. Pour $b \in \{1, \dots, B\}$ :

2.1 Pour  $i \in \{1, \dots, n\}$  : calculer les opposés du gradient :

$$U_i = - \frac{\partial}{\partial m(x_i)} \ell(m(x_i), y_i) \Big|_{m(x_i) = \hat{m}_{b-1}(x_i)} .$$

2.2 Estimer la règle faible sur l'échantillon  $((x_1, U_1), \dots, (x_n, U_n))$ .  
On obtient ainsi  $\hat{h}_b$ .

2.3 Mettre à jour :

$$\hat{m}_b(x) = \hat{m}_{b-1}(x) + \lambda \hat{h}_b(x)$$

où  $\lambda$  est un paramètre de régularisation.

## 3. Le prédicteur obtenu au final est :

$$\hat{m} = \hat{m}_B .$$

## Remarques

- On peut mettre en évidence l'agrégation en remarquant que :

$$\hat{g}_B(x) = \hat{g}_0(x) + \lambda \sum_{b=1}^B \hat{h}_b(x) .$$

- La vitesse de minimisation dépend des deux hyper-paramètres  $\lambda$  et  $B$  : si  $\lambda$  diminue,  $B$  augmente.



## Retour sur AdaBoost

- ▶ AdaBoost est équivalent à un problème de gradient boosting avec  $\lambda = 1$  et la fonction de perte :

$$\begin{aligned}\ell : \quad \mathbb{R} \times \{-1, 1\} &\rightarrow \mathbb{R}^+ \\ (u, y) &\mapsto \exp(-y u)\end{aligned}$$

- ▶  $\hat{g}(x)$  est un estimateur de :

$$g^*(x) = \frac{1}{2} \ln \left( \frac{\mathbb{P}(Y = 1 / X = x)}{\mathbb{P}(Y = -1 / X = x)} \right) .$$

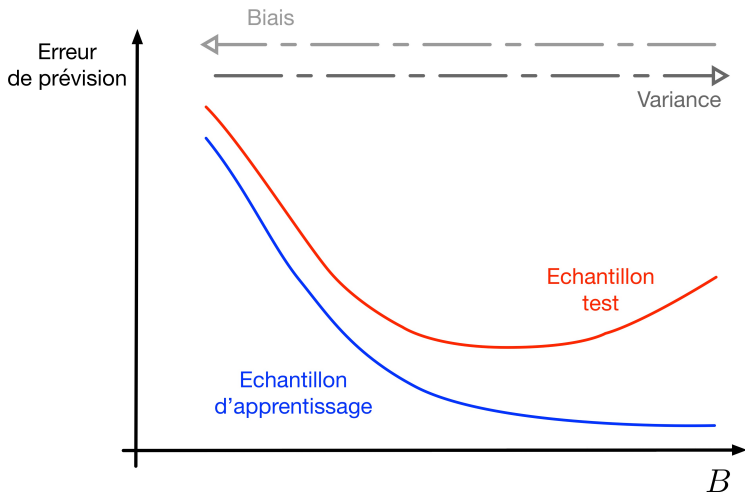
## Justification de la règle faible

- ▶ Le **boosting réduit le biais mais pas la variance**.
- ▶ La règle faible utilisée doit donc avoir une variance faible (et donc un biais élevé), comme les arbres avec très peu de feuilles.
- ▶ Considérer des règles non faibles n'améliore pas les performances de l'algorithme en général.

# Choix du nombre d'itérations I

- ▶ Le biais diminue avec  $B$ .
- ▶ Mais lorsque  $B$  est « trop » grand, la variance est importante : il existe un risque de sur-apprentissage.
- ▶ On choisit  $B$  par validation croisée.

## Choix du nombre d'itérations II



# Plan

Introduction

AdaBoost

Gradient boosting

Fonctions de pertes classiques

# LogitBoost I

- ▶ Le modèle de régression logistique, entre une variable  $Y \in \{0, 1\}$  et des covariables  $X = (X_1, \dots, X_p)^\top \in \mathbb{R}^p$ , pose :

$$p(x) := \mathbb{P}(Y = 1 / X = x) = \frac{\exp(\beta^\top x)}{1 + \exp(\beta^\top x)} = \frac{1}{1 + \exp(-\beta^\top x)}$$

où  $\beta \in \mathbb{R}^p$  est un paramètre inconnu estimé par maximum de vraisemblance.

- ▶ L'idée de **LogitBoost** est de supprimer l'hypothèse de linéarité de  $p(x)$ , en posant :

$$p(x) = \frac{1}{1 + \exp(-2g(x))}$$

où  $g : \mathbb{R}^p \rightarrow \mathbb{R}$  est une fonction inconnue.

## LogitBoost II

- La maximisation de :

$$\ln \left[ (p(x))^y (1 - p(x))^{1-y} \right] .$$

est équivalente à la minimisation de :

$$\ln [1 + \exp (-2\tilde{y}g(x))]$$

où  $\tilde{y} = 2y - 1 \in \{-1, 1\}$ .

- La fonction :

$$\ell : \mathbb{R} \times \{-1, 1\} \rightarrow \mathbb{R}^+$$

$$(u, y) \mapsto \ln (1 + \exp (-2yu))$$

est bien convexe en  $u$ .

## LogitBoost III

- ▶ LogitBoost désigne donc l'algorithme de gradient boosting avec la fonction de perte précédente.
- ▶ Comme pour AdaBoost, la solution  $\hat{g}(x)$  de LogitBoost est un estimateur de :

$$g^*(x) = \frac{1}{2} \ln \left( \frac{\mathbb{P}(Y = 1 / X = x)}{\mathbb{P}(Y = -1 / X = x)} \right) .$$

- ▶ Une fois la fonction  $\hat{g}$  déterminée, on obtient la règle de décision suivante :

$$y = \begin{cases} 1 & \text{si } \hat{p}(x) \geq \frac{1}{2} \\ -1 & \text{sinon} \end{cases} ,$$

soit

$$y = \begin{cases} 1 & \text{si } \hat{g}(x) \geq 0 \\ -1 & \text{sinon} \end{cases} .$$



## Autres exemples en classification supervisée

On peut considérer différentes fonctions  $\varphi$  telles que :

$$\ell(g(x), y) = \varphi(y g(x)) ,$$

par exemple :

- ▶ Fonction de perte quadratique :

$$\varphi(x) = (1 - x)^2 .$$

- ▶ Fonction de perte quadratique tronquée :

$$\varphi(x) = [(1 - x)_+]^2 .$$

- ▶ Fonction de perte Hinge (SVM) :

$$\varphi(x) = (1 - x)_+ .$$

## $L_2$ -Boosting I

- ▶ Le  $L_2$ -Boosting s'applique dans le cadre de la régression.
- ▶  $L_2$ -Boosting désigne l'algorithme de gradient boosting avec la fonction de perte :

$$\begin{aligned}\ell : \quad \mathbb{R}^2 &\rightarrow \mathbb{R}^+ \\ (u, y) &\mapsto \frac{1}{2} (u - y)^2\end{aligned}$$

- ▶ La solution  $\hat{m}(x)$  de  $L_2$ -Boosting est un estimateur de :

$$m^*(x) = \mathbb{E}(Y / X = x) .$$

## $L_2$ -Boosting II

- ▶ Dans l'étape 2.1 du gradient boosting, on obtient :

$$\begin{aligned} U_i &= -\frac{\partial}{\partial m(x_i)} \ell(m(x_i), y_i) \Big|_{m(x_i) = \hat{m}_{b-1}(x_i)} \\ &= y_i - \hat{m}_{b-1}(x_i) \end{aligned}$$

- ▶ Les  $U_i$  sont donc les résidus du prédicteur  $\hat{m}_{b-1}$  à l'étape  $(b-1)$ .
- ▶ Le prédicteur  $\hat{m}_b$  est élaboré via une régression sur les résidus de l'étape  $(b-1)$ .  
 $\hat{m}_{b-1}$  est amendée en expliquant l'information subsistant dans les résidus de l'étape  $(b-1)$ .

## Le coin R

On peut utiliser plusieurs packages, parmi lesquels :

- ▶ Le package `gbm` :
  - ▶ La fonction de base est `gbm`.
  - ▶ L'option `distribution` permet de spécifier la méthode : `adaboost` pour AdaBoost, `bernoulli` pour LogitBoost, `gaussian` pour  $L_2$ -Boosting, etc.
  - ▶ L'option `n.trees` permet de choisir le nombre d'arbres.
  - ▶ L'option `interaction.depth` permet de choisir la complexité des arbres.
  - ▶ L'option `shrinkage` correspond à  $\lambda$ .
  - ▶ La fonction `gbm.perf` permet de choisir le nombre d'itérations optimal.
- ▶ Le package `caret` :
  - ▶ La fonction de base est `train`.
  - ▶ On utilise une des options correspondant aux SVM, par exemple :
    - ▶ `method="ada"` pour AdaBoost,
    - ▶ `method="LogitBoost"` pour LogitBoost.

## Références

- Boyd, S. et L. Vandenberghe. 2003, *Convex optimization*, Cambridge University Press.
- Freund, Y. et R. E. Schapire. 1997, «A decision-theoretic generalization of on-line learning and an application to boosting», *Journal of Computer and System Sciences*, vol. 55, n° 1, p. 119–139.
- Schapire, R. E. et Y. Freund. 2012, *Boosting. Foundations and algorithms*, Adaptive Computation and Machine Learning, MIT Press.