



Travail d'étude et de recherche
réalisé dans le cadre du
Master 1 MAEF

RAPPORT DE TER

Boosting en apprentissage
statistique et **XGBoost**

Mohamed-Ayoub Bouzid
Alexis Sorre

Master 1 mathématiques appliquées
à l'économie et à la finance
à l'Université Paris 1 Panthéon Sorbonne

Encadrant : M. Jean-Marc Bardet
Février/Mai 2025

Table des matières

1	Introduction	3
2	Le modèle CART	5
2.1	Arbres de décision	5
2.2	Présentation du modèle CART	5
2.2.1	Arbres de régression	5
2.2.2	Arbres de classification	8
3	Théorie du Boosting	9
3.1	Motivation	9
3.2	Classificateur faible et théorème de Schapire	9
3.3	Schéma général de renforcement	14
4	AdaBoost	15
4.1	Présentation générale d'Adaboost	15
4.2	Description détaillée de l'algorithme	16
4.3	Analyse théorique	17
4.3.1	AdaBoost comme modèle additif	17
4.3.2	Fonction de perte exponentielle	18
4.3.3	Interprétation probabiliste	20
4.3.4	Limite de la perte exponentielle et problèmes de robustesse	21
5	XGBoost	23
5.1	Motivation et introduction	23
5.2	Principe général du Gradient Boosting	23
5.3	Du Gradient Boosting à XGBoost	25
5.4	XGBoost : Optimisation et Régularisation	26
5.4.1	Régularisation	26
5.4.2	Optimisation de la fonction objectif via l'approximation de Taylor	27
5.4.3	Critère de Split Optimal dans XGBoost	28
5.4.4	Parallélisation et Traitement des Valeurs Manquantes	28
5.5	Synthèse	29
5.6	Application numérique dans le cadre d'une régression	29
6	Simulation	33
6.1	Cadre de la simulation	33
6.2	Comportement modèle CART	33
6.3	Comportement d'AdaBoost	34
6.4	Comportement de XGBoost	35
6.5	Comparaison globale	35
7	Conclusion	37

1 Introduction

L'apprentissage statistique occupe une place centrale dans de nombreux domaines, allant de la médecine à la finance, en passant par l'industrie et les sciences sociales. Il s'agit d'un ensemble de méthodes permettant d'extraire des connaissances à partir de données, afin de faire des prédictions, de classer des observations, ou encore de détecter des relations cachées entre différentes variables. Les situations dans lesquelles ces techniques sont mobilisées sont extrêmement variées : estimer l'évolution future d'un cours boursier à partir d'indicateurs économiques, ou encore reconnaître automatiquement des chiffres manuscrits dans une image.

Dans tous ces cas, on dispose d'un ensemble d'observations appelées données d'apprentissage, contenant à la fois des résultats (variables à expliquer « targets ») et des caractéristiques associées (les variables explicatives « features »). Le but est alors de construire un modèle, ou apprenant (« learner »), capable de généraliser à de nouvelles situations en apprenant les régularités présentes dans ces données.

C'est précisément ce que l'on appelle l'apprentissage supervisé, une forme d'apprentissage statistique dans laquelle la variable cible guide l'entraînement du modèle.

Un problème d'apprentissage est dit supervisé lorsque les données d'apprentissage comprennent, pour chaque observation, les valeurs des variables explicatives ainsi que la valeur du résultat à prédire. L'apprentissage est donc guidé par cette variable cible, et le modèle apprend à établir un lien entre les entrées et la sortie attendue. Un exemple classique de problème d'apprentissage supervisé est la détection de courriels indésirables. Dans ce cas, chaque email représente une observation, et l'objectif est de prédire s'il s'agit d'un spam ou d'un message légitime. Pour cela, on dispose d'un ensemble d'emails préalablement étiquetés (on sait pour chacun s'il est un spam ou non (targets)), ainsi que dans notre cas la présence ou non de certains mots clés.

	george	you	your	hp	free	hpl	!	our	re	edu	remove
spam	0.00	2.26	1.38	0.02	0.52	0.01	0.51	0.51	0.13	0.01	0.28
email	1.27	1.27	0.44	0.90	0.07	0.43	0.11	0.18	0.42	0.29	0.01

FIGURE 1 – Ce tableau montre le pourcentage moyen d'apparition de certains mots ou caractères dans des messages spam par rapport à des emails normaux (non-spam). Les mots et caractères affichés sont ceux qui montrent la plus grande différence entre les deux types d'email

En analysant le tableau, on voit que le mot you apparaît en moyenne 2.26 fois dans un spam contre 1.27 dans un email.

Ainsi, notre méthode d'apprentissage doit décider de quel caractère utiliser et comment. Par exemple, on a :

```
if (%george < 0.4) & (%you > 1.5) then spam
else email
```

Selon la nature de la variable à prédire, on distingue deux grandes catégories de tâches supervisées :

1. La classification : Prédire une étiquette ou une classe parmi un ensemble fini (par exemple : spam / non spam, malade / sain).

2. La régression : Estimer une valeur numérique continue (comme le prix d'un bien, une température, ou une concentration chimique). Dans les deux cas, l'objectif est de développer un modèle prédictif performant, capable d'exploiter les régularités contenues dans les données passées pour effectuer des prédictions fiables sur de nouvelles données.

Une fois le cadre de l'apprentissage supervisé posé, il est essentiel de distinguer les différents types de modèles que l'on peut mobiliser pour cette tâche. On distingue notamment deux grandes approches : les modèles simples et les modèles d'ensemble.

Un modèle simple repose sur un unique algorithme d'apprentissage, appliqué directement aux données d'entraînement. C'est le cas, par exemple, de la régression linéaire, des arbres de décision. Ces modèles sont généralement faciles à mettre en œuvre, rapides à entraîner et souvent interprétables, ce qui en fait des outils privilégiés pour une première exploration des données ou pour des problèmes dont la structure reste relativement simple. Toutefois, leur capacité à capturer des

relations complexes ou des interactions non linéaires entre les variables peut être limitée, ce qui peut freiner leur performance sur certains jeux de données plus riches ou plus bruités. À l'inverse, les modèles d'ensemble reposent sur la combinaison de plusieurs modèles de base, souvent appelés apprenants faibles dans le but d'améliorer la robustesse des prédictions. L'idée est que, même si chaque modèle pris isolément est imparfait, leur agrégation permet de réduire les erreurs, de limiter le sur-apprentissage, et d'améliorer la généralisation sur de nouvelles données. Parmi les principales méthodes d'ensemble, on peut citer :

1. Le bagging, illustré notamment par les forêts aléatoires (**Random Forests**), qui consiste à entraîner plusieurs modèles sur des sous-échantillons aléatoires des données, puis à agréger leurs prédictions.

2. Le boosting, utilisé dans des algorithmes tels que **AdaBoost**, **Gradient Boosting** ou **XG-Boost**, qui corrige progressivement les erreurs des modèles précédents en les entraînant de manière séquentielle.

3. Le stacking, qui combine les prédictions de plusieurs modèles hétérogènes à l'aide d'un méta-modèle chargé d'apprendre comment les combiner de façon optimale.

Les modèles d'ensemble sont aujourd'hui largement utilisés en raison de leur excellente performance prédictive. Leur principal inconvénient réside toutefois dans leur complexité : ils sont généralement plus longs à entraîner, moins interprétables, et nécessitent une infrastructure de calcul plus robuste.

2 Le modèle CART

2.1 Arbres de décision

Avant de commencer à nous intéresser au modèle **CART** (**C**lassification and **R**egression **T**rees), il est essentiel de commencer par présenter ce que sont les arbres de décision.

Un arbre de décision est une méthode d'apprentissage supervisé dont l'idée principale est de partitionner l'espace des variables explicatives en plusieurs sous espaces et ajuster un modèle simple dans chacun d'eux (cf. figure 2). Par exemple, dans le cas d'une régression, on considère Y la variable endogène, X_1 et X_2 les variables explicatives et R_1, R_2, R_3 les sous espaces de la partition. Si $(X_1, X_2) \in R_n$ avec $n = 1, 2, 3$ alors on prédit Y par une constante c_m qui peut par exemple être la moyenne des Y dans cette région. Le modèle de régression associé sera alors de la forme :

$$\hat{f}(X) = \sum_{n=1}^3 c_n \mathbf{1}_{\{(X_1, X_2) \in R_n\}}.$$

Le partitionnement de cet espace peut être plus ou moins difficile à décrire toutefois dans le cadre de notre étude, nous nous limiterons aux partitions binaires récursives utilisés pour le modèle CART.

Ce type de partitionnement peut être représenté par un arbre binaire où à chaque noeud, nous avons deux branches, les données sont dirigées vers l'une des deux branches en fonction de la règle de décision choisie à chaque noeud. Chaque chemin de l'arbre résulte en une feuille donnant une prédiction finale.

Le principal avantage d'utiliser d'utiliser les arbres binaires récursifs est leur interprétabilité. En effet, on peut décrire la partition sous la forme d'un seul arbre tandis que lorsqu'on a plus de variables d'entrées la représentation peut être bien plus difficile.

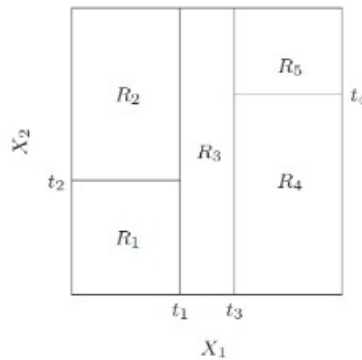


FIGURE 2 – Partitionnement de l'espace des variables explicatives constitué de X_1 et X_2

2.2 Présentation du modèle CART

Le modèle **CART** est une méthode d'apprentissage supervisé utilisée pour des tâches de classification et de régression. Il permet d'expliquer comment une variable peut être expliquée en fonction d'autres facteurs. Nous allons présenter ici les deux catégories d'arbres de décision à savoir les arbres de régression et de classification.

2.2.1 Arbres de régression

On va considérer dans cette partie Y une variable à prédire, p variables endogènes X_1, \dots, X_p et N observations (x_i, y_i) avec $i = 1, \dots, N$. Une première méthode que l'on peut envisager pour trouver la meilleure prédiction de Y possible est de choisir la partition qui minimise le risque empirique. Si on suppose que l'espace des variables se décompose en 2 partitions R_1, R_2 pour simplifier et que la prédiction de Y est c_m dans la région m avec $m \in \{1, 2\}$. Le modèle est donc de la forme :

$$f(x) = c_1 \mathbf{1}_{\{x \in R_1\}} + c_2 \mathbf{1}_{\{x \in R_2\}}$$

Nous voulons alors minimiser :

$$\text{SSE}(R_1, R_2) = \sum_{x_i \in R_1} (y_i - f(x_i))^2 + \sum_{x_i \in R_2} (y_i - f(x_i))^2.$$

D'autre part, nous déterminons les estimateurs \hat{c}_m des c_m avec $m = 1, 2$ en minimisant :

$$\sum_{x_i \in R_m} (y_i - f(x_i))^2$$

On considère que $\forall x_i \in R_m$ on a $f(x_i) = c$. On a donc $Q_m(c) = \sum_{x_i \in R_m} (y_i - c)^2$. On va trouver le minimum en résolvant :

$$\frac{dQ_m(c)}{dc} = 0$$

En développant, on a :

$$Q_m(c) = \sum_{x_i \in R_m} (y_i - c)^2 = \sum_{x_i \in R_m} (y_i^2 - 2y_i c + c^2) = \sum_{x_i \in R_m} y_i^2 + 2c \sum_{x_i \in R_m} y_i + c^2 N_m$$

où $N_m = \text{card}(x_i \in R_m)$. Donc en dérivant l'expression on obtient :

$$\frac{dQ_m(c)}{dc} = -2 \sum_{x_i \in R_m} y_i + 2c N_m = 0$$

puis en isolant c , on trouve finalement que la valeur optimale de c est la moyenne des y_i dans la région R_m :

$$\hat{c}_m = \frac{1}{N_m} \sum_{x_i \in R_m} y_i \quad (1)$$

On vérifie qu'il s'agit bien d'un minimum en montrant que Q_m est bien convexe :

$$\frac{d^2 Q_m(c)}{dc^2} = 2N_m > 0.$$

Cependant, il apparaît que cette méthode est irréalisable car trop coûteuse en calculs. En effet, elle implique de tester toutes les façons possibles de découper l'espace i.e 2^N possibilité dans le cas où on veut découper l'espace en 2 partitions.

Détermination des variables et seuils de séparation :

Une solution à ce problème réside dans l'utilisation du modèle **CART**. L'idée derrière cet algorithme est que l'on cherche à prédire Y en déterminant des variables de séparation X_j , les points de séparation et la forme que l'arbre doit avoir. A chaque étape, l'algorithme choisit la variable et le point de séparation qui minimisent un critère d'erreur. Ce processus est répété jusqu'à atteindre un critère d'arrêt.

On considère donc deux demi plans :

$$R_1(j, s) = \{X \mid X_j < s\} \quad \text{et} \quad R_2(j, s) = \{X \mid X_j > s\}. \quad (2)$$

Pour une étape donnée, on va chercher la variable de séparation j et le point de séparation s qui résolvent :

$$\min_{j, s} \left[\min_{c_1} \sum_{x_i \in R_1(j, s)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in R_2(j, s)} (y_i - c_2)^2 \right] \quad (3)$$

Commençons par minimiser dans un premier temps la fonction :

$$L(c_1) = \sum_{x_i \in R_1(j, s)} (y_i - c_1)^2$$

Cette fonction étant convexe, on obtient en minimum on trouve un minimum en résolvant l'équation :

$$\frac{d}{dc_1} \sum_{x_i \in R_1(j,s)} (y_i - c_1)^2 = 0$$

Après développement puis dérivation, on obtient :

$$\sum_{x_i \in R_1(j,s)} y_i + 2c_1 \sum_{x_i \in R_1(j,s)} 1 = 0$$

Finalement, on obtient :

$$c_1 = \frac{1}{|R_1|} \sum_{x_i \in R_1(j,s)} y_i \quad (4)$$

Ce qui correspond à la moyenne des y_i dans R_1 :

$$\hat{c}_1 = \text{moyenne}(y_i | x_i \in R_1(j,s)) \quad (5)$$

On va également minimiser :

$$L(c_2) = \sum_{x_i \in R_2(j,s)} (y_i - c_2)^2$$

Par un raisonnement analogue, on trouve également :

$$\hat{c}_2 = \text{moyenne}(y_i | x_i \in R_2(j,s)) \quad (6)$$

Ces valeurs \hat{c}_1 et \hat{c}_2 sont calculés pour chaque couple (j,s) puis on évalue l'erreur quadratique :

$$\sum_{x_i \in R_1(j,s)} (y_i - \hat{c}_1)^2 + \sum_{x_i \in R_2(j,s)} (y_i - \hat{c}_2)^2$$

On choisit le couple (j,s) qui donne l'erreur quadratique la plus petite qui seront utilisés pour réaliser la première division de l'arbre. Le processus est ensuite répété sur toutes les régions résultantes.

Détermination de la taille optimale de l'arbre :

Le deuxième élément cruciale dans la construction d'un arbre de régression est la détermination de sa taille. En effet, une fois que l'arbre qu'on notera T_0 a été construit à l'aide de la méthode précédemment décrit, il se peut que cet arbre soit trop grand ce qui aura comme conséquence un surajustement du modèle : l'arbre sera très performant sur les données d'entraînement mais mal généralisable aux nouvelles données car trop spécifique. Pour remédier à cela, nous allons tenter d'obtenir une sorte d'optimalité dans la taille de l'arbre.

Une première idée serait de diviser l'arbre en un noeud que si la somme des carrés des résidus en ce noeud ne dépasse pas un certain seuil que nous fixons. Cependant, cette approche ne semble pas pertinente car il se peut que pour un noeud donné la somme des carrés des résidus dépasse ce seuil mais que en continuant de diviser on obtient des prédictions rendant cette somme des carrés des résidus.

Nous allons au lieu de cela élaguer l'arbre avec la méthode "élagage par complexité-coût". On considère $T \subset T_0$ un arbre obtenu par élagage à partir de T_0 . On indexe chaque noeud terminal par m . Nous définissons, $|T|$ comme le nombre de noeuds terminaux dans T , $N_m = \text{card}(x_i \in R_m)$ et \hat{c}_m comme dans l'équation. De plus, l'erreur quadratique moyenne sur R_m est donnée par :

$$Q_m(T) = \frac{1}{N_m} \sum_{x_i \in R_m} (y_i - \hat{c}_m)^2 \quad (7)$$

Le **critère de complexité-coût** est défini par :

$$C_\alpha(T) = \sum_{m=1}^{|T|} N_m Q_m(T) + \alpha |T| \quad (8)$$

Pour une valeur de $\alpha \geq 0$, l'arbre optimal T_α sera celui qui minimise ce critère :

$$T_\alpha = \arg \min_T C_\alpha(T) \quad (9)$$

Nous devons choisir un α réalisant un compromis entre la taille de l'arbre et sa qualité d'ajustement. Si α est grand, les grands arbres sont pénalisés et T_α sera petit. Si α est petit, l'ajustement est privilégié et plus de noeuds sont gardés.

2.2.2 Arbres de classification

Dans de nombreuses situations, la variable à prédire Y n'est pas une variable numériques mais qualitatifs : malade ou non, client fidèle ou à risque, spam ou message légitime. Ces problèmes relèvent tous de la classification. Nous considérons que Y prend les valeurs $1, \dots, K$. Ces valeurs représentent les différentes classes considérées.

Le principe est ici le même que précédemment c'est à dire que nous commençons par construire le grand arbre T_0 dont le critère d'arrêt pour sa construction par exemple si le nombre minimal d'observation dans un noeud est atteint. Nous verrons qu'un autre critère d'arrêt peut également être utilisé. Pour chaque variable d'entrée X_1, \dots, X_p , on va tester tous les seuils de coupure possibles. Les seuils de coupures sont déterminés en commençant par trier les observations pour une variable explicative donnée et ensuite nous faisons la moyenne entre deux observations consécutives. Si on N observations pour chaque variable explicative, nous obtenons $N - 1$ seuils de coupure à tester. Nous introduisons désormais de nouvelles notions. Nous considérons pour commencer un noeud m associé à une région R_m contenant N_m observations. Nous définissons **la proportion d'observation de classe k dans le noeud m** comme :

$$\hat{p}_{mk} = \frac{1}{N_m} \sum_{x_i \in R_m} 1(y_i = k) \quad (10)$$

Un noeud sera dit impur s'il ne contient pas uniquement des observations d'une seule classe. On définit donc **la classe majoritaire dans le noeud m** comme :

$$k(m) = \arg \max_k \hat{p}_{mk} \quad (11)$$

Pour chaque noeud, l'objectif sera de mesurer l'impureté et de trouver le couple de variables de coupure et de seuil de coupure qui minimisent l'impureté globale (définie ci-dessous). Les deux mesures d'impureté communément utilisées sont l'indice Gini et l'entropie croisée. Ils sont respectivement définis par :

$$\text{Indice Gini} = \sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk}) \quad (12)$$

et

$$\text{Entropie croisée} = - \sum_{k=1}^K \hat{p}_{mk} \log(\hat{p}_{mk}) \quad (13)$$

On utilise l'une de ces deux mesures pour mesurer l'impureté globale d'une division et on choisira la division qui la minimise :

$$\text{Impureté Globale} = \frac{N_L}{N} Q_L + \frac{N_R}{N} Q_R \quad (14)$$

avec N_L, N_R respectivement le nombre d'observations dans les sous-noeuds gauche et droit et N le nombre totale d'observation dans les deux sous noeuds.

On répète récursivement ce processus dans chaque sous noeud jusqu'à atteindre un critère d'arrêt. Nous en avons évoqué un ci-dessus. Il est à noter toutefois qu'un autre critère pouvant être utilisé dans le cas de classification est le gain d'impureté c'est à dire qu'on s'arrête si une division ne diminue pas de manière assez significative l'impureté globale.

Une fois cet arbre, construit nous allons l'élaguer comme précédemment c'est à dire en trouvant l'arbre T minimisant le critère complexité-coût mais cette fois $Q_m(T)$ correspondra à la mesure d'impureté dans le noeud m .

3 Théorie du Boosting

3.1 Motivation

Les arbres de décision pour la régression et la classification permettent une prédiction assez intuitives de la variable cible Y . Toutefois pris isolément, il apparaît que ces modèles de prédiction présentent de nombreuses limites. L'une des principales limitations est que ces modèles sont sensibles aux variations des données. Une modification, même minime, de l'échantillon d'apprentissage peut engendrer une structure significativement différente. Cela peut avoir pour conséquence une variance élevée des prédictions compromettant ainsi la robustesse des prédictions et limitant la capacité de généralisation du modèle.

Par ailleurs, lorsque les arbres de décision sont profonds, le nombre de partition possible explose de manière exponentielle. Pour un prédicteur ayant q valeurs possibles le nombre de partitions possibles est de 2^q pouvant provoquer des difficultés de computation un surapprentissage du modèle. Présentons maintenant une méthode permettant de travailler non pas avec un seul arbre complexe mais combiner plusieurs arbres peu profond afin d'obtenir un apprenant global plus stable et souvent plus robuste.

3.2 Classificateur faible et théorème de Schapire

Dans cette partie, nous nous intéressons à ce qui constitue la base du boosting à savoir les classificateurs faibles. Ce sont ces classificateurs qui sont combinés pour former un apprenant global plus performant. Pour montrer que cela est possible, nous allons énoncer le **théorème de Schapire** et en proposer une esquisse de preuve.

Définition 1. *Un classificateur faible est un algorithme d'apprentissage qui donne des résultats au moins meilleurs que le hasard. Cela signifie que sa précision est supérieure à 50% dans un problème de classification binaire. Les classificateurs faibles considérés sont souvent de petits arbres CART peu profond.*

Afin d'énoncer ce qu'est le théorème de Schapire, nous avons besoin de poser le cadre théorique nécessaire. On se place dans le cadre du modèle PAC introduit par Valiant en 1984. Dans ce cadre, un **concept** c est une fonction borélienne $c : X \rightarrow \{0, 1\}$ attribuant à chaque donnée une étiquette. De plus, une **classe de concept** C est un ensemble de ces fonctions. On peut la décomposer en sous-classe C_n de sorte que chacune de ces sous-classes agit sur un ensemble de données X_n . Par ailleurs, on considère que l'algorithme n'a pas accès directement à la fonction c mais uniquement à l'ensemble d'entraînement. De plus, l'algorithme reçoit des données $x \in X_n$ choisit aléatoirement selon une distribution inconnue D , et pour chaque donnée, la réponse correcte $c(x)$ est connue et déterminée par la fonction cible. Cela permet de traduire le fait qu'en pratique, l'algorithme apprend à partir d'un échantillon aléatoire sans connaître la loi qui régit la génération des données.

A partir des données observées, l'algorithme va produire une fonction $h : X_n \rightarrow \{0, 1\}$ appelée hypothèse. C'est cette fonction qui sert à prédire les étiquettes de nouvelles données. Nous mesurerons la qualité de cette hypothèse avec le taux d'erreur :

$$\text{err}_D(h) = P(h(x) \neq c(x)) \quad (15)$$

Lorsque cette erreur est inférieure à $\epsilon > 0$, on dit que h est ϵ -proche du concept cible. Nous définissons de façon plus formelle ce que sont les apprenants forts et faibles.

Définition 2. *Une classe de concept C est dite fortement apprenable s'il existe un algorithme capable pour tout concept $c \in C$, toute distribution D , et pour tout seuil d'erreur ϵ et de confiance δ , de produire avec une probabilité d'au moins $1 - \delta$ une hypothèse h telle que :*

$$\text{err}_D(h) \leq \epsilon$$

Définition 3. *Un algorithme est dit faiblement apprenable s'il peut produire une hypothèse légèrement meilleur que le hasard :*

$$\text{err}_D(h) \leq \frac{1}{2} - \gamma$$

Enfin, le dernier concept clé dans ce cadre est celui de **l'oracle d'exemple EX**.

Définition 4. *Un oracle d'exemple pour un concept c sur un domaine X_n , associé à une distribution D , est une procédure qui à chaque appel tire une instance $x \in X_n$ selon une distribution D et retourne la paire $(x, c(x))$, i.e l'instance et son étiquette correcte donnée par le concept cible c . L'oracle d'exemple est donc ce qui génère les données d'entraînement.*

Le cadre théorique étant posé, nous pouvons maintenant énoncer le théorème de Schapire :

Théorème 1. *Toute classe de concepts qui est faiblement apprenable est également fortement apprenable. Plus précisément, si un algorithme d'apprentissage produit une hypothèse dont l'erreur est inférieure à $\frac{1}{2} - \gamma$ pour une constante $\gamma > 0$ alors il est possible de construire à partir de celui-ci, un nouvel algorithme qui produit une hypothèse avec une erreur arbitrairement petite.*

Preuve. L'idée générale de la preuve que propose Schapire est qu'il construit un algorithme appelé **Learn** qui simule un algorithme faible A trois fois sur des distributions modifiées des données d'entraînement.

```

Learn( $\epsilon, \delta, EX$ )
  Input:   error parameter  $\epsilon$ 
            confidence parameter  $\delta$ 
            examples oracle  $EX$ 
            (implicit) size parameters  $s$  and  $n$ 

  Return: hypothesis  $h$  that is  $\epsilon$ -close to the target concept  $c$  with probability  $\geq 1 - \delta$ 

  Procedure:
    if  $\epsilon \geq 1/2 - 1/p(n, s)$  then return WeakLearn( $\delta, EX$ )
     $\alpha \leftarrow g^{-1}(\epsilon)$ 

     $EX_1 \leftarrow EX$ 
     $h_1 \leftarrow \text{Learn}(\alpha, \delta/5, EX_1)$ 
     $\tau_1 \leftarrow \epsilon/3$ 
    let  $\hat{a}_1$  be an estimate of  $a_1 = \Pr_{v \in D}[h_1(v) \neq c(v)]$ :
      choose a sample sufficiently large that  $|a_1 - \hat{a}_1| \leq \tau_1$  with probability  $\geq 1 - \delta/5$ 
    if  $\hat{a}_1 \leq \epsilon - \tau_1$  then return  $h_1$ 

    defun  $EX_2()$ 
      { flip coin
        if heads, return the first instance  $v$  from  $EX$  for which  $h_1(v) = c(v)$ 
        else return the first instance  $v$  from  $EX$  for which  $h_1(v) \neq c(v)$  }
     $h_2 \leftarrow \text{Learn}(\alpha, \delta/5, EX_2)$ 
     $\tau_2 \leftarrow (1 - 2\alpha)\epsilon/8$ 
    let  $\hat{e}$  be an estimate of  $e = \Pr_{v \in D}[h_2(v) \neq c(v)]$ :
      choose a sample sufficiently large that  $|e - \hat{e}| \leq \tau_2$  with probability  $\geq 1 - \delta/5$ 
    if  $\hat{e} \leq \epsilon - \tau_2$  then return  $h_2$ 

    defun  $EX_3()$ 
      { return the first instance  $v$  from  $EX$  for which  $h_1(v) \neq h_2(v)$  }
     $h_3 \leftarrow \text{Learn}(\alpha, \delta/5, EX_3)$ 

    defun  $h(v)$ 
      {  $b_1 \leftarrow h_1(v)$ ,  $b_2 \leftarrow h_2(v)$ 
        if  $b_1 = b_2$  then return  $b_1$ 
        else return  $h_3(v)$  }
    return  $h$ 

```

FIGURE 3 – Algorithme Learn extrait de "The strength of weak learnability" (Schapire,1990), p.7

Nous commençons par fixer l'erreur α tel que $g(\alpha) = \epsilon$ avec $g(\alpha) = 3\alpha^2 - 2\alpha^3$. On appelle une première fois $\text{Learn}(\alpha, \frac{\delta}{5}, EX_1)$ et on considère un estimateur \hat{a}_1 de $a_1 = P(h_1(v) \neq c(v))$ qui soit suffisamment proche de a_1 avec une probabilité supérieure ou égale à $1 - \frac{\delta}{5}$. Si, $\hat{a}_1 \leq \epsilon$ alors on renvoie h_1 car cela signifierait que l'algorithme est un apprenant fort. Sinon, on construit une distribution

EX₂ qui donne avec 50% de chance un exemple bien classé par h_1 et 50% un exemple mal classé. Cela va forcer le prochain appel à se concentrer autant sur les erreurs de h_1 que sur ses succès. On appelle ensuite $\text{Learn}(\alpha, \frac{\delta}{5}, \text{EX}_2)$ pour obtenir une hypothèse h_2 et si elle est suffisamment bonne on la retourne. Sinon, on construit une troisième distribution EX₃ qui sélectionne uniquement les exemples où $h_1 \neq h_2$ puis on entraîne une troisième hypothèse avec $\text{Learn}(\alpha, \frac{\delta}{5}, \text{EX}_3)$. Enfin, si $h_1 = h_2$ alors on choisit h_1 sinon on choisit h_3 .

N.B : Nous remarquons que lorsque Learn est appelé récursivement, $\frac{\delta}{5}$ est pris en paramètre au lieu de δ . Cela est dû au fait qu'il y a 5 sources d'échecs lors de l'exécution de l'algorithme : l'appel récursif pour h_1 , l'estimation statistique de l'erreur \hat{a}_1 , l'appel récursif de h_2 , l'estimateur de \hat{e} et l'appel récursif pour h_3 . Si chaque évènement a une probabilité inférieure ou égale à $\frac{\delta}{5}$ alors $P(\text{au moins une étape échoue}) \leq 5 \cdot \frac{\delta}{5} = \delta$. Donc la probabilité que l'algorithme réussisse est supérieur à $1 - \delta$.

Nous montrons maintenant le résultat suivant permettant de conclure la preuve du théorème 1 :

Théorème 2. Soit $0 < \epsilon < \frac{1}{2}$ et $0 \leq \delta \leq 1$. Alors l'algorithme $\text{Learn}(\epsilon, \delta, EX)$ retourne avec une probabilité au moins $1 - \delta$, une hypothèse h telle que $P(h(x) \neq c(x)) \leq \epsilon$

Démonstration. Nous avons déjà expliqué dans la remarque pourquoi la probabilité qu'une bonne exécution ait lieu est au moins $1 - \delta$. Il nous reste à montrer que si Learn effectue une bonne exécution, alors l'hypothèse de sortie est ϵ -proche du concept cible.

1^{er} cas :

Dans le cas où l'estimateur \hat{a}_1 ou \hat{e} sont trouvés plus petit que $\epsilon - \tau_1$ ou $\epsilon - \tau_2$ respectivement avec \hat{a}_1 et \hat{e} suffisamment proche de a_1 et e respectivement alors l'hypothèse retournée est ϵ -proche du concept cible.

Cas général :

Soit a_i l'erreur de h_i sous la distribution D_i induite par l'oracle EX _{i} au $i^{\text{ème}}$ appel récursif. Par hypothèse de récurrence nous avons pour $i = 1, 2, 3$, $a_i \leq \alpha$ car Learn est appelé récursivement avec un paramètre d'erreur α . On introduit les notations suivantes :

$$p_i(v) = P(h_i(v) \neq c(v)) \quad (16)$$

$$q(v) = P(h_1(v) \neq h_2(v)) \quad (17)$$

$$w = \sum_{v \sim D} P(h_2(v) \neq h_1(v) = c) \quad (18)$$

$$x = \sum_{v \sim D} P(h_1(v) = h_2(v) = c) \quad (19)$$

$$y = \sum_{v \sim D} P(h_1(v) \neq h_2(v) = c) \quad (20)$$

$$z = \sum_{v \sim D} P(h_1(v) = h_2(v) \neq c) \quad (21)$$

Remarque et notation : La notation $\sum_{v \sim D} P(\pi(v))$ désigne la probabilité qu'un prédicteur τ soit vérifié sur des instances v tirées de X_n selon la distribution D .

X_n est l'ensemble des instances possibles pour le problème d'apprentissage.

Puis par la formule des probabilités totales, nous obtenons :

$$\sum_{v \sim D} P(\tau(v)) = \sum_{v \in X_n} P(v)P(\tau(v)|v) = \sum_{v \in X_n} D(v)P(\tau(v)) \quad (22)$$

Dans la suite, afin d'alléger la notation, P_D représentera $\sum_{v \sim D} P$.

Nous exprimons maintenant les lois de probabilité selon lesquelles les oracles EX _{i} génèrent les instances v . Cela nous permettra ensuite de calculer l'erreur finale de l'hypothèse h construite par Learn .

Soit $D(v)$ la probabilité que $v \in X_n$ soit tirés aléatoirement l'oracle d'exemple EX alors on a que :

$$D(v) = D_1(v) \quad (23)$$

Pour la distribution D_2 , l'oracle d'exemple EX_2 est construit en lançant une pièce équilibrée, puis si la pièce tombe sur pile, on retourne la première instance v tel que $h_1(v) = c(v)$. Sinon, l'algorithme retourne la première instance v de EX telle que $h_1(v) \neq c(v)$. Nous en déduisons donc que :

$$D_2(v) = \frac{1}{2}P(v|h_1(v) = c(v)) + \frac{1}{2}P(v|h_1(v) \neq c(v)) \quad (24)$$

Or par la formule de Bayes, nous avons :

$$P(v|h_1(v) \neq c(v)) = \frac{P(v \text{ et } h_1(v) \neq c(v))}{P(h_1(v) \neq c(v))} = \frac{D(v)p_1(v)}{a_1} \quad (25)$$

$$P(v|h_1(v) = c(v)) = \frac{D(v)(1 - p_1(v))}{1 - a_1} \quad (26)$$

avec $1 - a_1 = w + x = P_D(h_1(v) = c(v))$.

Finalement, en injectant ces expressions dans (23), on obtient :

$$D_2(v) = \frac{D(v)}{2} \left(\frac{p_1(v)}{a_1} + \frac{1 - p_1(v)}{1 - a_1} \right) \quad (27)$$

La distribution D_3 ne considère que les instances où les hypothèses h_1 et h_2 sont en désaccord. Nous avons donc :

$$D_3(v) = P(v|h_1(v) \neq h_2(v)) = \frac{P(v \text{ et } h_1(v) \neq h_2(v))}{P(h_1(v) \neq h_2(v))} = \frac{D(v)q(v)}{w + y} \quad (28)$$

Nous calculons maintenant l'erreur de l'hypothèse h construite à partir de h_1, h_2, h_3 i.e

$$P_D(h(v) \neq c(v)) \quad (29)$$

où $c(v)$ est la bonne étiquette de v .

D'après la définition de h , deux cas peuvent se produire :

1er cas : $h_1(v) = h_2(v)$, donc $h(v) = h_1(v) = h_2(v)$. L'algorithme fait une erreur si cette prédiction commune est fausse.

Deuxième cas : $h_1(v) \neq h_2(v)$ et dans ce cas $h(v) = h_3(v)$. Dans ce cas, l'algorithme fait une erreur si $h_3(v) \neq c(v)$.

On obtient donc :

$$P_D(h(v) \neq c(v)) = P_D(h_1(v) = h_2(v) = h(v)) + P_D((h_1(v) \neq h_2(v)) \cap (h_3(v) \neq c(v))) \quad (30)$$

Puis en utilisant le résultat de la remarque, nous avons :

$$P_D(h(v) \neq c(v)) = z + \sum_{v \in X_n} D(v)q(v)p_3(v) = z + \sum_{v \in X_n} (w + y)D_3(v)q(v)p_3(v) \quad (31)$$

Le lemme suivant donne les résultats qui permettent de conclure la preuve :

Lemme 1. *Nous avons les 3 équations suivantes :*

$$1 - a_2 = \frac{y}{2a_1} + \frac{z}{2(1 - a_1)} \quad (32)$$

$$z = a_1 - y \quad (33)$$

$$w = (2a_2 - 1)(1 - a_1) + \frac{y(1 - a_1)}{a_1} \quad (34)$$

Démonstration. On commence par montrer (32). Nous avons par définition :

$$1 - a_2 = P_{D_2}(h_2(v) = c(v)) = \sum_{v \in X_n} D_2(v)(1 - p_2(v))$$

Puis en injectant l'expression de $D_2(v)$ dans la somme, nous obtenons

$$\begin{aligned}
1 - a_2 &= P_{D_2}(h_2(v) = c(v)) = \sum_{v \in X_n} \left(\frac{1}{2} \frac{D(v)p_1(v)}{a_1} + \frac{1}{2} \frac{D(v)(1-p_1(v))}{1-a_1} \right) (1-p_2(v)) \\
&\iff 1 - a_2 = \sum_{v \in X_n} \frac{1}{2} \frac{D(v)p_1(v)}{a_1} (1-p_2(v)) + \frac{1}{2} \frac{D(v)(1-p_1(v))}{1-a_1} (1-p_2(v))
\end{aligned}$$

Nous obtenons finalement en utilisant encore une fois la remarque que :

$$1 - a_2 = \frac{y}{2a_1} + \frac{z}{2(1-a_1)}$$

Nous démontrons maintenant (33). Il est clair que :

$$\{h_1(v) \neq h_2(v), h_2(v) = c(v)\} \cap \{h_1(v) = h_2(v), h_2(v) = c(v)\} = \emptyset$$

Nous avons donc :

$$\begin{aligned}
y + z &= P_D(\{h_1(v) \neq h_2(v), h_2(v) = c(v)\}) + P_D(\{h_1(v) = h_2(v), h_2(v) = c(v)\}) \\
&= P_D(\{h_1(v) \neq h_2(v), h_2(v) = c(v)\} \cup \{h_1(v) = h_2(v), h_2(v) = c(v)\}) \\
&= P_D(\{h_1(v) \neq c(v)\}) = a_1
\end{aligned}$$

Nous terminons la démonstration de ce lemme en montrant (34). On sait que par définition, on a :

$$\begin{aligned}
w &= \sum_{v \in X_n} D(v)(1-p_1(v))p_2(v) \\
&= \sum_{v \in X_n} D(v)(1-p_1(v))(1-(1-p_2(v))) \\
&= (1-a_1) - \sum_{v \in X_n} D(v)(1-p_1(v))(1-p_2(v)) \\
&= (1-a_1) - z
\end{aligned}$$

D'autre part, en reprenant (32), nous avons :

$$\begin{aligned}
1 - a_2 &= \frac{y}{2a_1} + \frac{z}{2(1-a_1)} \\
\iff 2(1-a_2) &= \frac{y}{a_1} + \frac{z}{1-a_1} \\
\iff z &= (1-a_1) \left(2(1-a_2) - \frac{y}{a_1} \right)
\end{aligned}$$

Nous obtenons donc :

$$\begin{aligned}
w &= (1-a_1) - (1-a_1) \left(2(1-a_2) - \frac{y}{a_1} \right) \\
&= (1-a_1) \left(1 - \left(2(1-a_2) - \frac{y}{a_1} \right) \right) \\
&= (1-a_1) \left(2a_1 - 1 + \frac{y}{a_1} \right) \\
&= (2a_2 - 1)(1-a_1) + \frac{y(1-a_1)}{a_1}
\end{aligned}$$

□

Nous pouvons maintenant conclure la preuve de ce théorème en repartant de l'équation (31), nous obtenons :

$$\begin{aligned}
P_D(h(v) \neq c(v)) &= z + a_3(w + y) \\
&\leq z + \alpha(w + y) \quad \text{par récursivité} \\
&= a_1 - y + \alpha(2a_2 - 1)(1 - a_1) + \alpha \frac{y(1 - a_1)}{a_1} + \alpha y \quad \text{en injectant (33) et (34)} \\
&= \alpha(2a_2 - 1)(1 - a_1) + a_1 + y \left(-1 + \alpha \left(1 + \frac{1 - a_1}{a_1} \right) \right) \\
&= \alpha(2a_2 - 1)(1 - a_1) + a_1 + \frac{y(\alpha - a_1)}{a_1} \\
&= \alpha(2a_2 - 1)(1 - a_1) + a_1 + y \left(-1 + \alpha \left(1 + \frac{1 - a_1}{a_1} \right) \right) \\
&\leq \alpha(2a_2 - 1)(1 - a_1) + a_1 + \alpha - a_1 \quad \text{car } y \leq a_1 \\
&= \alpha(2\alpha - 1)(1 - \alpha) \\
&= 3\alpha^2 - 2\alpha^3 = g(\alpha) = \epsilon
\end{aligned}$$

Ce qui conclut la preuve du théorème 2. \square

D'après le théorème 2, nous savons que l'algorithme $\text{Learn}(\epsilon, \delta, \text{EX})$ retourne avec une probabilité au moins $1 - \delta$ une hypothèse h tel que $\text{err}_D(h) \leq \epsilon$. Cela nous permet d'établir que tout algorithme faible peut être transformé en algorithme fort, c'est à dire capable d'atteindre une erreur arbitrairement petite. \square

3.3 Schéma général de renforcement

Maintenant que nous avons expliqué ce que sont les classificateurs faibles et établi le **théorème de Schapire**, nous proposons dans cette section une approche plus intuitive et algorithmique du boosting.

Le boosting repose sur l'idée de construire un classificateur faible qui peut par exemple être construit à partir du modèle CART, que l'on combine pour obtenir un modèle global plus performant. A chaque itération, un nouveau modèle est entraîné pour corriger les erreurs du précédents de sorte à ce que chaque étape vient affiner les prédictions du modèle précédent. Nous pouvons alors décrire le boosting comme une construction additive du type :

$$f(x) = \sum_{m=1}^M \beta_m b(x, \gamma_m) \quad (35)$$

avec $b(x, \gamma_m)$ sont les apprenants faibles et β_m sont les coefficients déterminant leur importance dans le modèle final. Cette construction est appelé **forward stagewise additive modelling** et peut être vu comme une somme de corrections successives apportées à un modèle initialement trivial.

L'intuition centrale du boosting est que l'attention du modèle est déplacée vers les observations mal prédites. Autrement dit, à chaque étape, on entraîne un classificateur sur une version pondérée des données où les observations bien prédites voient leur poids diminuer tandis que celles mal prédites voient leurs poids augmenter.

Le boosting modifie donc implicitement à chaque itération la distribution des données en accentuant l'importance des cas difficiles. Cette idée déjà présente dans la construction théorique du boosting de Schapire en 1990 se retrouve également au coeur des algorithmes modernes tel que **AdaBoost** ou **XGBoost**.

Dans la suite, nous allons présenter plus en détails les deux plus importants algorithmes de boosting que sont AdaBoost et XGBoost.

4 AdaBoost

4.1 Présentation générale d'Adaboost

AdaBoost est un algorithme de boosting introduit en 1997 par Yoav Freund et Robert Schapire. Il a été l'un des premiers algorithmes qui a traduit les fondements théoriques du boosting en une méthode concrète.

Dans AdaBoost, les modèles faibles les plus fréquemment utilisés sont des arbres de décision de profondeur 1, appelés stumps. Ce sont donc des arbres qui effectuent une seule coupure sur une variable explicative ce qui les rend très simple mais également très limités individuellement. Nous les construisons, en utilisant la méthode de construction des arbres de classification décrite dans la partie traitant du modèle **CART** c'est à dire que le stump cherche la variable et le seuil qui mesure d'impureté comme l'indice Gini ou l'entropie croisée. Néanmoins, puisqu'il ne peut utiliser qu'une seule variable, un stump ne peut capter qu'une information partielle sur le problème.

Exemple : Si nous souhaitons prédire le risque de défaut d'un emprunteur, nous pourrions utiliser plusieurs caractéristiques telles que le revenu, la situation professionnelle, le statut marital etc... Un stump ne peut se baser que sur l'une de ces variables ce qui limite fortement sa prédiction. Ce genre d'arbre est considéré comme des classificateurs faibles. AdaBoost entraîne de façon séquentielle ces stumps en modifiant de façon dynamique le poids des observations à chaque itération. Au départ, chaque observation a le même poids et un premier stump noté G_1 est entraîné sur ses données pondérées. Ensuite, les poids sont mis à jour de façon à ce que les observations mal classées voient leurs poids augmenter tandis que les observations bien classées voient leur poids diminuer. Le stump suivant se concentrera davantage sur les observations avec un poids plus important et ce processus est répété un nombre de fois M fixé. Ainsi, l'ordre des stumps a une importance car chaque modèle dépend de ceux qui le précèdent. A la fin, les stumps sont combinés via un vote pondéré où chaque stump reçoit un poids α_m . Dans le cas d'une classification binaire la prédiction finale est donnée par :

$$G(x) = \text{sign} \left(\sum_{m=1}^M \alpha_m G_m(x) \right) \quad (36)$$

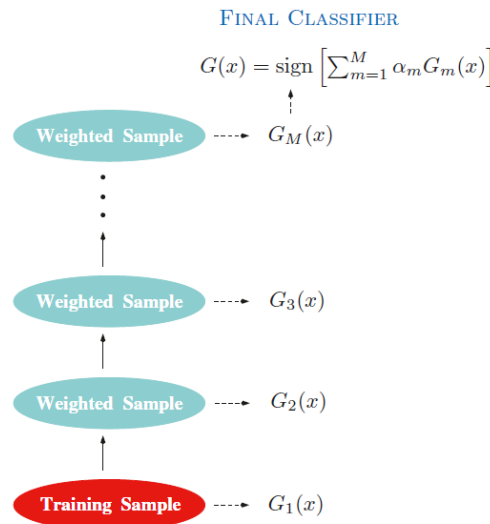


FIGURE 4 – Schéma de fonctionnement d'AdaBoost

Source : Hastie, Tibshirani, Friedman – *The Elements of Statistical Learning*, 2^e édition, Springer (2009), Figure 10.1.

4.2 Description détaillée de l'algorithme

Nous présentons dans cette partie l'algorithme **AdaBoost** et nous donnons également un exemple illustratif permettant de bien le comprendre. Nous considérons un problème de classification à deux classes. La variable à prédire Y prend donc deux valeurs : 1 et -1. Pour un "feature" X , le classificateur produit $G(X)$ à valeurs dans $\{-1, 1\}$. L'algorithme AdaBoost se présente comme suit :

1. Initialize the observation weights $w_i = 1/N$, $i = 1, 2, \dots, N$.
2. For $m = 1$ to M :
 - (a) Fit a classifier $G_m(x)$ to the training data using weights w_i .
 - (b) Compute

$$\text{err}_m = \frac{\sum_{i=1}^N w_i I(y_i \neq G_m(x_i))}{\sum_{i=1}^N w_i}.$$
 - (c) Compute $\alpha_m = \log((1 - \text{err}_m)/\text{err}_m)$.
 - (d) Set $w_i \leftarrow w_i \cdot \exp[\alpha_m \cdot I(y_i \neq G_m(x_i))]$, $i = 1, 2, \dots, N$.
3. Output $G(x) = \text{sign} \left[\sum_{m=1}^M \alpha_m G_m(x) \right]$.

FIGURE 5 – Pseudo code de l'algorithme AdaBoost.M1.

Source : Hastie, Tibshirani, Friedman – *The Elements of Statistical Learning*, 2^e édition, Springer (2009), Algorithm 10.1.

Nous considérons les observations $\{(x_1, y_1), \dots, (x_N, y_N)\}$. Nous commençons dans un premier temps à attribuer à toutes les observations un poids identique :

$$w_i^{(1)} = \frac{1}{N} \text{ pour } i = 1, \dots, N$$

Ensuite, pour l'itération $m \in \{1, \dots, M\}$, nous entraînons le classificateur faible $G_m(x)$ sur les données d'entraînement pondérées selon les poids $w_i^{(m)}$. Les erreurs sur les observations les plus pondérées sont pénalisées davantage. Pour cette itération, nous calculons l'erreur pondérée qui est définie par :

$$\text{err}_m = \frac{\sum_{i=1}^N w_i I(y_i \neq G_m(x_i))}{\sum_{i=1}^N w_i} \quad (37)$$

L'intérêt de normaliser l'erreur de prédiction est d'avoir une interprétation probabiliste de cette erreur qui ne peut prendre que des valeurs comprises dans l'intervalle $[0, 1]$. Cela nous permet donc d'avoir une interprétation homogène de la performance du modèle courant. Autrement dit, cela permet de faire en sorte que l'erreur reste à la même échelle et donc que les poids conservent une cohérence à chaque itération.

Une fois que nous calculons l'erreur pondérée, err_m du classificateur G_m , nous quantifions son importance dans le modèle final. Nous calculons son poids α_m définie par :

$$\alpha_m = \log \left(\frac{1 - \text{err}_m}{\text{err}_m} \right) \quad (38)$$

Ce poids va permettre de mesurer la fiabilité du classificateur. En effet, plus err_m sera faible, plus α_m sera grand augmentant ainsi son influence dans la prédiction finale. A l'inverse, lorsque err_m se rapproche de 0.5 c'est à dire que le classificateur a une précision équivalente à celle du hasard alors α_m tend vers 0. Il est important de noter que la forme de α_m n'est pas arbitraire et découle de la minimisation d'une fonction de perte exponentielle dans un modèle additif. Nous justifions en détail cette affirmation dans la section 4.3.

Nous mettons ensuite à jour le poids des observations w_i . L'idée à cette étape est d'augmenter l'importance des observations mal classées et de réduire l'importance de celles qui ont été bien classées. En effet, si une observation i est bien classée alors $I(y_i \neq G_m(x)) = 0$ et

$$w_i \exp(0) = w_i$$

donc son poids reste inchangé. A l'inverse si une observation est mal classée alors $I(y_i \neq G_m(x)) = 1$ et son poids est multiplié par un facteur supérieur ou égale à 1. Intuitivement, cela signifie que l'algorithme va davantage sur les observations mal classées. Nous expliquerons la forme de la mise à jour des poids w_i dans la section 4.3. Enfin, après M itérations, le classificateur final sera :

$$G(x) = \text{sign} \left(\sum_{m=1}^M \alpha_m G_m(x) \right) \quad (39)$$

Chaque G_m vote pour la classe 1 ou -1 sur la donnée x et son influence est déterminée par son poids α_m . Les votes des classificateurs les plus performants comptent davantage dans le vote final.

Remarque : Dans l'algorithme présenté, nous utilisons deux notions de poids qu'il est important de ne pas confondre. D'une part, w_i représente le poids de l'observation i dans l'apprentissage du prochain classificateur. Il est ajusté de manière à ce que les observations mal classées obtiennent un poids plus important. Par ailleurs, α_m est le poids du classificateur m et détermine son influence dans la prédiction finale.

Exemple : On considère le jeu de données suivant contenant trois observations :

x_i	y_i
A	1
B	-1
C	1

FIGURE 6 – Jeu de données simplifié à trois observations pour illustrer le fonctionnement de l'algorithme AdaBoost

Nous appliquons AdaBoost pendant 2 itérations. A la première itération, les trois observations ont un poids $w_i^{(1)} = \frac{1}{3}$ pour $i = 1, 2, 3$. Nous supposons que le premier classifieur G_1 prédit correctement les points 1 et 3 mais se trompe sur le point 2. En calculant l'erreur pondérée on trouve $\text{err}_m = \frac{1}{3}$ et le poids du classifieur est $\alpha_1 = \log(2)$. Enfin actualisant les poids w_i , nous avons $w_1 = \frac{1}{3}, w_2 = \frac{2}{3}, w_3 = \frac{1}{3}$. Donc nous retrouvons bien l'idée que les observations mal classées se voient attribués un poids plus important. A la deuxième itération, nous supposons que le classifieur G_2 prédit correctement toutes les observations. Nous avons $\text{err}_2 = 0$ et $\alpha_2 \rightarrow +\infty$. En pratique, nous fixons un ϵ très petit par exemple $\epsilon = 10^{-10}$ donc $\alpha_2 \approx 10 \log(10) \approx 23$. Finalement, la sortie de l'algorithme sera $G(x) = \text{sign}(\alpha_1 G_1(x) + \alpha_2 G_2(x))$ et comme $\alpha_2 \gg \alpha_1$, le vote du classificateur faible aura plus d'importance en cohérence avec ce qui a été expliqué.

4.3 Analyse théorique

4.3.1 AdaBoost comme modèle additif

Un modèle additif est un modèle qui s'écrit sous la forme :

$$f(x) = \sum_{m=1}^M \beta_m b(x, \gamma_m) \quad (40)$$

avec β_m les coefficients d'expansion et $b(x, \gamma_m)$ est une fonction définie sur x et paramétré par γ_m . Il apparaît donc que dans ce modèle, les paramètres inconnus sont les β_m et γ_m . Afin de les estimer,

nous voudrions déterminer la fonction f qui minimise une fonction de perte sur l'échantillon d'apprentissage c'est à dire qui permettrait d'obtenir la plus grande précision sur la prédiction globale. Le problème d'optimisation serait donc le suivant :

$$\min_{\{(\beta_m, \gamma_m)\}_{m=1}^M} \sum_{i=1}^N L \left(y_i, \sum_{m=1}^M \beta_m b(x_i; \gamma_m) \right) \quad (41)$$

Résoudre ce problème d'optimisation nous permettrait d'ajuster tous les paramètres simultanément de sorte à minimiser l'erreur globale. Toutefois, il apparait que ce problème est souvent infaisable en terme de complexité de calcul. Une alternative va consister à déterminer les paramètres (β_m, γ_m) en minimisant la fonction de perte par étapes successives. Cette procédure est connue sous le nom de **modélisation additive par étapes**. Elle est décrite dans l'algorithme suivant qui constitue une méthode générale d'optimisation pour les modèles additifs.

1. Initialize $f_0(x) = 0$.

2. For $m = 1$ to M :

(a) Compute

$$(\beta_m, \gamma_m) = \arg \min_{\beta, \gamma} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + \beta b(x_i; \gamma)).$$

(b) Set $f_m(x) = f_{m-1}(x) + \beta_m b(x; \gamma_m)$.

FIGURE 7 – Pseudo-code de la modélisation additive par étapes

Source : Hastie, Tibshirani, Friedman – *The Elements of Statistical Learning*, 2^e édition, Springer (2009), Algorithm 10.2.

Nous commençons par initialiser la fonction f à 0 c'est à dire que nous ne faisons aucune prédiction au départ. Puis pour chaque étape $m \in \{1, \dots, M\}$, nous recherchons parmi toutes les fonctions de base possible $b(x, \gamma)$ et tous les coefficients β permettant de réduire la perte globale si nous l'ajoutons au modèle. Autrement dit, nous cherchons la fonction de base qui permet d'améliorer le plus la prédiction globale. Ceci se traduit dans l'algorithme par :

$$(\beta_m, \gamma_m) = \arg \min_{\beta, \gamma} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + \beta b(x_i; \gamma)) \quad (42)$$

Une fois ce couple de paramètres trouvés, nous mettons à jour le modèle :

$$f_m(x) = f_{m-1}(x) + \beta_m b(x, \gamma_m) \quad (43)$$

A la fin des M itérations, nous trouvons la forme finale de f .

A partir de l'équation (39), nous pouvons voir qu'**AdaBoost** peut être interprété comme un modèle additif construit de manière incrémentale. Dans notre cas, les fonctions $b(x, \gamma_m)$ correspondent aux classificateurs faibles $G_m(x) \in \{-1, 1\}$ et les coefficients β_m sont les poids α_m attribués à chaque classifieurs. Nous allons montrer dans la partie qui suit, qu'AdaBoost correspond à la minimisation de la fonction de perte exponentielle au sein de ce cadre additif.

4.3.2 Fonction de perte exponentielle

Nous allons montrer dans cette partie qu'AdaBoost est équivalent à l'algorithme de modélisation additive par étapes lorsqu'on choisit la fonction de perte :

$$L(x, f(x)) = \exp(-yf(x)) \quad (44)$$

qu'on appelle **fonction de perte exponentielle**.

Démonstration. Nous reprenons l'algorithme présenté à la figure 6 et nous supposons qu'à l'itération $m - 1$, nous avons construit le modèle partiel $f_{m-1}(x)$. Nous cherchons alors le terme $\beta G(x)$ qui permet d'obtenir $f_m(x) = f_{m-1}(x) + \beta G(x)$ de sorte à minimiser la perte totale :

$$(\beta_m, G_m) = \arg \min_{\beta, G} \sum_{i=1}^N \exp(-y_i (f_{m-1}(x_i) + \beta G(x_i))) \quad (45)$$

Nous définissons les poids mesurant à quel point l'observation i est mal classée par :

$$w_i^{(m)} = \exp(-y_i f_{m-1}(x_i)) \quad (46)$$

Lorsque la prédiction est incorrecte, $y_i f_{m-1}(x_i) < 0$ et par croissance de la fonction exponentielle le poids $w_i^{(m)}$ est plus élevé que lorsque la prédiction est correcte en cohérence avec ce que nous avons expliqué lorsque nous avons traité l'algorithme **AdaBoost**. On réécrit alors (45) de la façon suivante :

$$(\beta_m, G_m) = \arg \min_{\beta, G} \sum_{i=1}^N w_i^{(m)} \exp(-y_i \beta G(x_i)) \quad (47)$$

Par ailleurs, comme on considère que $y_i, G(x_i) \in \{-1, 1\}$, on a alors que $\exp(-\beta y_i G(x_i)) = \exp(-\beta)$ si $y_i = G(x_i)$ et $\exp(-\beta y_i G(x_i)) = \exp(\beta)$ si $y_i \neq G(x_i)$. Nous posons :

$$\begin{aligned} C &= \{i : y_i = G(x_i)\} \\ M &= \{i : y_i \neq G(x_i)\} \end{aligned}$$

L'équation (47) devient alors :

$$(\beta_m, G_m) = \arg \min_{\beta, G} \exp(-\beta) \sum_{i \in C} w_i^{(m)} + \exp(\beta) \sum_{i \in M} w_i^{(m)} \quad (48)$$

Par ailleurs, on sait aussi que :

$$\sum_{i=1}^N w_i^{(m)} = \sum_{i \in C} w_i^{(m)} + \sum_{i \in M} w_i^{(m)} \quad (49)$$

Puis en injectant l'expression de $\sum_{i \in C} w_i^{(m)}$ dans (48), on obtient :

$$\begin{aligned} \exp(-\beta) \sum_{i \in C} w_i^{(m)} + \exp(\beta) \sum_{i \in M} w_i^{(m)} &= \exp(-\beta) \left(\sum_{i=1}^N w_i^{(m)} - \sum_{i \in M} w_i^{(m)} \right) + \exp(\beta) \sum_{i \in M} w_i^{(m)} \\ &= \exp(-\beta) \sum_{i=1}^N w_i^{(m)} + (\exp(\beta) - \exp(-\beta)) \sum_{i \in M} w_i^{(m)} \end{aligned}$$

Pour β fixé, minimiser (48) revient à minimiser la somme pondérée des erreurs :

$$G_m = \arg \min_G \sum_{i \in M} w_i^{(m)} = \arg \min_G \sum_{i=1}^N w_i^{(m)} I(y_i \neq G(x_i)) \quad (50)$$

ce qui est exactement fait par AdaBoost i.e à chaque itération G_m est ajusté de sorte à minimiser la somme pondérée des erreurs. Nous fixons maintenant G_m et nous définissons :

$$\text{err}_m = \frac{\sum_{i=1}^N w_i^{(m)} I(y_i \neq G_m(x_i))}{\sum_{i=1}^N w_i^{(m)}} \quad (51)$$

En reprenant encore une fois l'expression à minimiser dans (48), nous obtenons :

$$L(\beta) = \exp(-\beta)(1 - \text{err}_m) \sum_{i=1}^N w_i^{(m)} + \exp(\beta) \text{err}_m \sum_{i=1}^N w_i^{(m)} \quad (52)$$

Finalement pour obtenir la valeur de β qui minimise L , nous résolvons l'équation :

$$\begin{aligned} \frac{dL}{d\beta}(\beta) = 0 &\iff -(1 - \text{err}_m) \exp(-\beta) + \text{err}_m \exp(\beta) = 0 \\ &\iff \exp(-2\beta) = \frac{\text{err}_m}{1 - \text{err}_m} \\ &\iff \beta_m = \frac{1}{2} \log \left(\frac{1 - \text{err}_m}{\text{err}_m} \right) \end{aligned}$$

C'est exactement la formule à un facteur près du poids du classificateur m utilisée dans AdaBoost. Finalement, pour une étape donnée, l'algorithme se termine par la mise à jour du modèle :

$$f_m(x) = f_{m-1}(x) + \beta_m G_m(x) \quad (53)$$

Le poids à l'étape $m + 1$ devient donc :

$$w_i^{(m+1)} = \exp(-y_i f_m(x_i)) = w_i^{(m)} \exp(-\beta_m y_i G_m(x_i)) \quad (54)$$

et en utilisant le fait que $-y_i G_m(x_i) = 2I(y_i \neq G_m(x_i)) - 1$, nous obtenons donc que :

$$w_i^{(m+1)} = w_i^{(m)} \exp(\alpha_m \mathbb{I}(y_i \neq G_m(x_i))) \exp(-\beta_m) \quad (55)$$

avec $\alpha_m = 2\beta_m$. Comme le facteur $\exp(-\beta_m)$ multiplie tous les poids, il n'a aucun effet. \square

Nous venons donc de montrer que la modélisation additive par étape avec une fonction de perte exponentielle est équivalente à appliquer l'algorithme AdaBoost. Nous allons maintenant essayer de comprendre pour la fonction de perte exponentielle permet de retrouver Adaboost et en tirer une interprétation probabiliste.

4.3.3 Interprétation probabiliste

Dans la section précédente, nous avons présenté la perte exponentielle empirique c'est à dire que les différentes minimisations s'effectuent sur un jeu de données $\{(x_i, y_i)\}_{i=1}^N$. Nous allons maintenant, dans cette section considérer la forme théorique de la perte exponentielle afin de déterminer la fonction idéale qu'AdaBoost essaie d'estimer. Nous introduisons alors la fonction de perte exponentielle théorique comme :

$$L(f) = \mathbb{E}_{Y|x}[\exp(-Yf(x))] \quad (56)$$

où l'espérance est prise selon la loi de Y conditionnellement à l'entrée x . Cela est justifié par le fait que nous connaissons l'entrée x et nous voulons chercher la valeur optimale de $f^*(x)$ qui minimise la perte moyenne sur la distribution de Y sachant cette entrée x . Nous allons alors chercher la fonction f^* tel que :

$$f^*(x) = \arg \min_f \mathbb{E}_{Y|x}[\exp(-Yf(x))] \quad (57)$$

Comme $Y \in \{-1, 1\}$, nous notons $p = P(Y = 1|x)$ et donc $1 - p = P(Y = -1|x)$. Nous obtenons donc :

$$L(f) = p \cdot \exp(-f) + (1 - p) \cdot \exp(f)$$

Pour minimiser cette expression, nous résolvons :

$$\begin{aligned} \frac{dL}{df} &= -p \cdot \exp(-f) + (1 - p) \cdot \exp(f) = 0 \\ &\iff (1 - p) \cdot \exp(f) = p \cdot \exp(-f) \\ &\iff \exp(2f) = \frac{p}{1 - p} \\ &\iff f^*(x) = \frac{1}{2} \log \left(\frac{p}{1 - p} \right) \end{aligned}$$

Ce résultat peut se réécrire comme :

$$f^*(x) = \frac{1}{2} \log \left(\frac{P(Y = 1|x)}{P(Y = -1|x)} \right) \quad (58)$$

ou de façon équivalente :

$$\Pr(Y = 1 | x) = \frac{1}{1 + e^{-2f^*(x)}} \quad (59)$$

Nous constatons alors que le score f appris par AdaBoost est la moitié d'un logit de la probabilité conditionnelle $P(Y = 1|x)$. Donc Adaboost peut également être interprété comme un estimateur probabiliste apprenant une approximation du logit.

Par ailleurs, nous introduisons un autre critère utilisé pour modéliser la probabilité conditionnelle $P(Y = 1|x)$ est la log-vraisemblance binomiale négative également appelée déviance définie par :

$$l(Y, f(x)) = \log(1 + \exp(-2Yf(x))) \quad (60)$$

Cette fonction va nous permettre de réaliser une comparaison de robustesse avec la perte exponentielle car nous pouvons montrer que cette fonction est minimisée comme la perte exponentielle par :

$$f^*(x) = \frac{1}{2} \log \left(\frac{P(Y = 1|x)}{P(Y = -1|x)} \right)$$

et donc visent le même objectif théorique.

4.3.4 Limite de la perte exponentielle et problèmes de robustesse

Nous nous intéressons dans cette section à la robustesse d'**Adaboost** et nous allons notamment la comparer avec celle de la déviance. Nous rappelons que la fonction score $f(x)$ associée à une observation permet de mesurer la confiance du modèle dans sa prédiction. Plus cette fonction score est grande et plus le modèle a confiance. Lorsque $yf(x) > 0$, l'observation est bien classée et lorsque cette quantité est négative elle est mal classée. Quand $yf(x) = 0$, le modèle se trouve exactement sur la frontière de décision. Par ailleurs, nous savons que l'objectif de tout algorithme de classification est de produire des marges positives aussi fréquemment que possible.

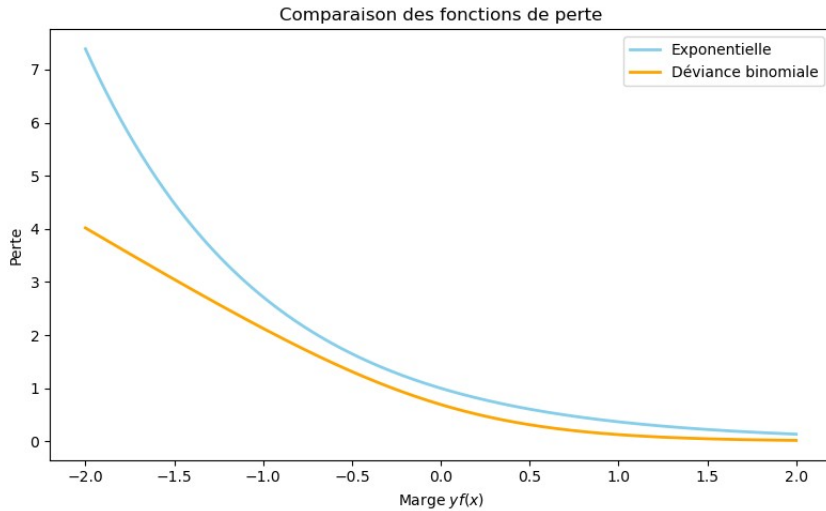


FIGURE 8 – Comparaison entre la perte exponentielle et la déviance binomiale en fonction de la marge $yf(x)$.

La figure ci-dessus compare la perte en fonction de la marge $yf(x)$ de la fonction de perte exponentielle et de la déviance. Nous pouvons observer que la déviance augmente de façon assez linéaire pour les marges négatives tandis que la perte exponentielle croît de façon exponentielle pour ces

même marges. Cela implique que le critère exponentielle concentre une influence plus forte sur les observations mal classées que la déviance qui elle, répartit l'influence plus uniformément sur toutes les observations.

Cela rend donc la déviance nettement plus robuste dans un contexte bruité. En effet, une observation est mal étiquetée ou aberrante alors la perte exponentielle va donner un poids significatif à cette observation quitte à modifier fortement le modèle. Cela peut donc détériorer les prédictions sur tout le reste du jeu de données. La perte exponentielle possède alors un risque de surapprentissage localisé. Cela est moins le cas de la déviance qui prend en compte ce genre de points mais sans leur données une importance significative.

Nous venons ici de mettre en évidence certaines limites d'AdaBoost et notamment sa grande sensibilité aux observations mal classées. Nous allons alors chercher à développer un algorithme plus robuste et donc nous intéresser à XGBoost.

5 XGBoost

5.1 Motivation et introduction

Après avoir présenté **AdaBoost**, qui repose sur une pondération adaptative des observations mal classées pour combiner plusieurs classificateurs faibles, nous allons introduire une approche plus générale et plus puissante : **le Gradient Boosting**.

L'idée fondamentale derrière le Gradient Boosting est de reformuler le problème du boosting comme un problème d'optimisation, où chaque nouvel apprenant faible est construit dans le but de réduire l'erreur résiduelle du modèle précédent, à l'aide de la descente de gradient.

Plutôt que de se focaliser uniquement sur les erreurs de classification comme dans AdaBoost, le Gradient Boosting adopte une approche plus flexible : il minimise une fonction de perte différentiable quelconque (par exemple, l'erreur quadratique en régression) en ajoutant itérativement des prédicteurs qui pointent dans la direction du négatif du gradient de la fonction de perte, calculé par rapport aux prédictions actuelles du modèle.

C'est dans ce cadre qu'est introduit **XGBoost** (eXtreme Gradient Boosting) par Tianqi Chen et Carlos Guestrin en 2016 dans *XGBoost : A Scalable Tree Boosting System*, une version optimisée et hautement performante du Gradient Boosting, largement utilisée en pratique et ayant gagné de nombreuses compétitions de science des données.

Comparé au Gradient Boosting classique et à d'autres méthodes basées sur les arbres, XGBoost présente deux avantages majeurs.

Premièrement, il est efficace sur le plan computationnel, car il traite les données sous forme de blocs compressés, ce qui permet un tri et un traitement rapides en parallèle. Cela lui confère un temps d'entraînement nettement inférieur à celui du Gradient Boosting classique (cf. graphique ci-dessous). Deuxièmement, il utilise un développement de Taylor d'ordre deux pour minimiser la fonction objectif. Cela signifie qu'il prend en compte la courbure de la fonction objectif, ce qui lui permet de converger vers un minimum de meilleure qualité que d'autres méthodes.

Il est également important de noter que XGBoost ne s'appuie pas sur les critères classiques utilisés pour la construction des arbres de décision, tels que l'entropie croisée ou l'indice de Gini définis précédemment. Il utilise à la place des critères spécifiques, qui sont présentés et expliqués en détail dans ce document.

Troisièmement, XGBoost permet via des mécanismes de régularisation intégrés de limiter le surapprentissage.

Enfin, il permet également de gérer automatiquement les données manquantes de manière intelligente. Lors de l'entraînement, le modèle apprend quelle direction (gauche ou droite dans l'arbre) est la plus optimale pour une valeur manquante, ce qui évite d'avoir à forcer une imputation artificielle.

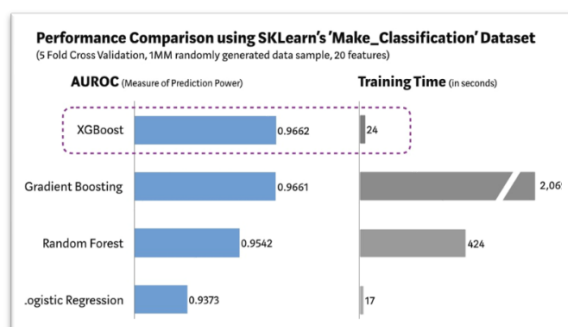


FIGURE 9 – Enter Caption

5.2 Principe général du Gradient Boosting

Le Gradient Boosting repose sur l'utilisation séquentielle de plusieurs arbres de décision. En effet, contrairement à la forêt aléatoire, les arbres de décision ne sont pas construits indépendamment. Le

nouvel arbre est entraîné par les erreurs commises par l'ensemble des arbres précédents. Pour chaque observation, chaque arbre fournit une prédiction, et la prédiction finale est obtenue en additionnant les contributions de l'ensemble des arbres.

Avant de discuter du Gradient Boosting, il est important de revenir sur l'idée du boosting glouton. Le boosting traditionnel consiste à ajouter des arbres de manière séquentielle, avec l'objectif de corriger les erreurs faites par les arbres précédents. À chaque itération, un nouvel arbre est ajouté pour améliorer la prédiction globale. Cette méthode peut être considérée comme gloutonne dans la mesure où, à chaque étape, l'algorithme se concentre uniquement sur l'erreur locale (les erreurs résiduelles des arbres précédents) sans considérer de manière globale la minimisation de la fonction de perte.

Le Gradient Boosting améliore cette approche en introduisant une optimisation plus globale qui prend en compte la fonction de perte dans son ensemble, en suivant le gradient de celle-ci. L'idée est de guider l'algorithme dans une direction plus précise à chaque itération pour minimiser la fonction de perte de manière optimale, et non simplement de réduire les erreurs résiduelles locales.

Le Gradient Boosting repose sur un concept fondamental : suivre le gradient de la fonction de perte. Cela signifie que chaque nouvel arbre construit à l'itération m est ajusté non seulement pour corriger les erreurs résiduelles du modèle précédent, mais pour suivre la direction du gradient de la fonction de perte. Le gradient représente la direction dans laquelle la fonction de perte diminue le plus rapidement, et c'est dans cette direction que l'arbre doit être ajusté. Pour calculer ce gradient, on prend la dérivée de la fonction de perte par rapport à la prédiction du modèle, ce qui nous donne une idée claire de l'ajustement nécessaire pour réduire l'erreur.

Au début du processus de boosting, chaque nouvel arbre est ajouté pour corriger les erreurs des prédictions précédentes. À chaque itération m , on calcule le gradient de la fonction de perte $L(y_i, f(x_i))$ par rapport aux prédictions actuelles $f_{m-1}(x_i)$ (après $m-1$ itérations), c'est-à-dire celles issues des arbres précédemment ajoutés. Ce gradient est défini comme suit :

$$g_{im} = - \left. \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right|_{f(x_i)=f_{m-1}(x_i)}$$

Cela correspond à la dérivée de la fonction de perte par rapport à la prédiction actuelle du modèle. Le gradient indique la direction dans laquelle les prédictions doivent être ajustées pour minimiser l'erreur. Plus précisément, il montre à quel point il faut ajuster les prédictions pour réduire la fonction de perte. C'est dans cette direction que le nouvel arbre sera construit.

Une fois le gradient calculé, un nouvel arbre de décision $T_m(x_i)$ est construit. Cet arbre a pour objectif de minimiser l'erreur sur les résidus, c'est-à-dire d'ajuster les prédictions en suivant la direction du gradient. L'ajustement de l'arbre permet de corriger les erreurs du modèle actuel en suivant la pente négative du gradient de la fonction de perte.

L'arbre est ajusté de manière à réduire les erreurs sur les données, et la prédiction est mise à jour comme suit :

$$f_m(x) = f_{m-1}(x) + T_m(x)$$

Le modèle est mis à jour après l'ajout de chaque arbre. La prédiction totale $f_M(x)$ après M itérations est donnée par la somme des prédictions successives :

$$f_M(x) = \sum_{m=1}^M T(x; \Theta_m)$$

À chaque itération, chaque arbre $T(x; \Theta_m)$ est ajusté de manière à minimiser les erreurs des arbres précédents, en suivant la direction du gradient de la fonction de perte. Cela permet de réduire l'erreur globale de manière plus efficace que dans le boosting traditionnel.

Le calcul des régions R_j et des prédictions γ_j est au cœur de l'optimisation du boosting par gradient. Chaque région R_j est déterminée pour minimiser la fonction de perte $L(y_i, \gamma_j)$, en ajustant les partitions de l'espace de manière à réduire les résidus.

Les régions R_j peuvent être calculées de manière itérative, et la solution optimale pour chaque γ_j dans une région R_j est trouvée par un algorithme de partitionnement de type régression, qui vise à minimiser l'erreur quadratique dans chaque région.

Dans le gradient boosting, une approche de **descente de gradient** est utilisée pour résoudre le problème d'optimisation. La descente de gradient est une méthode itérative qui permet de trouver les paramètres Θ_m qui minimisent la fonction de perte. À chaque itération, le modèle est ajusté en suivant la direction du gradient de la fonction de perte.

La mise à jour du modèle $f_m(x)$ en ajoutant un arbre $T_m(x)$ est réalisée en utilisant la solution de l'optimisation suivante :

$$\hat{\Theta}_m = \arg \min_{\Theta_m} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + T(x_i; \Theta_m))$$

Cela permet d'ajuster progressivement le modèle pour réduire l'erreur sur l'ensemble des données, en ajoutant des arbres successifs qui suivent la direction du gradient.

On a ainsi l'algorithme de Boosting d'Arbres de Décision par Gradient :

1. **Initialisation** : Initialiser le modèle $f_0(x) = \arg \min_{\gamma} \sum_{i=1}^N L(y_i, \gamma)$, avec une prédiction simple (par exemple, la moyenne des y_i).
2. **Pour** $m = 1$ à M :
 - (a) Pour $i = 1, 2, \dots, N$, calculer les résidus r_{im} comme suit :

$$r_{im} = - \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x_i)=f_{m-1}(x_i)}.$$

- (b) Ajuster un arbre de régression $T(x; \Theta_m)$ aux cibles r_{im} , ce qui donne les régions terminales $R_{jm}, j = 1, 2, \dots, J_m$.
- (c) Pour $j = 1, 2, \dots, J_m$, calculer les prédictions γ_{jm} en résolvant :

$$\gamma_{jm} = \arg \min_{\gamma} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \gamma).$$

- (d) Mettre à jour la prédiction du modèle :

$$f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm}).$$

3. **Sortie** : La prédiction finale du modèle est donnée par $\hat{f}(x) = f_M(x)$.

5.3 Du Gradient Boosting à XGBoost

On a vu que le Gradient Boosting repose sur l'idée de construire des modèles successifs pour réduire les erreurs résiduelles des arbres précédents, en ajustant les prédictions vers la direction du gradient de la fonction de perte. Cela est réalisé en suivant une approche itérative où, à chaque itération, un nouvel arbre est ajouté pour améliorer les prédictions. Bien que

ce processus soit puissante, il présente certaines limites. Il peut être lent à l'exécution sur de grands volumes de données en raison de son apprentissage séquentiel. Ce modèle est également sujet à un risque de surapprentissage si les arbres sont trop complexes. De plus, il gère mal les valeurs manquantes, nécessitant souvent une imputation préalable. Enfin, son optimisation repose uniquement sur le gradient de premier ordre, sans prendre en compte la courbure de la fonction de perte, ce qui peut limiter sa précision.

XGBoost (eXtreme Gradient Boosting) surmonte ces limitations en introduisant plusieurs améliorations notables, tant sur le plan de l'optimisation que de la régularisation.

5.4 XGBoost : Optimisation et Régularisation

XGBoost prend les principes du Gradient Boosting et les améliore avec plusieurs optimisations clés, qui sont détaillées ci-dessous.

5.4.1 Régularisation

L'une des principales innovations de XGBoost par rapport au Gradient Boosting classique est l'ajout de régularisation dans la fonction objectif. La régularisation permet de mieux contrôler la complexité des arbres et d'éviter le surajustement, en limitant la profondeur des arbres et la norme des poids.

La fonction objectif dans XGBoost comprend donc un terme de régularisation, $\Omega(T)$, ajouté à la fonction de perte. La fonction objectif totale devient ainsi :

$$\mathcal{L}(\Theta) = \sum_{i=1}^N L(y_i, \hat{y}_i) + \lambda \sum_{j=1}^J \|\gamma_j\|^2 + \gamma \sum_{k=1}^K d_k$$

- $\sum_{i=1}^N L(y_i, \hat{y}_i)$: fonction de perte qui mesure l'écart entre la prédiction \hat{y}_i et la vérité y_i . Cette fonction doit être différentiable et convexe (afin de pouvoir être optimisée facilement et atteindre ainsi son minimum)

- Pénalisation des poids :

$$\lambda \sum_{j=1}^J \|\gamma_j\|^2$$

Cette somme pénalise la norme des poids (ou scores) associés aux feuilles des arbres. Cela évite que les poids prennent des valeurs trop extrêmes.

- Signification des symboles :

- γ_j : poids attribué à la feuille j dans un arbre.
- $\|\gamma_j\|^2$: carré de la norme (souvent γ_j^2 si γ_j est scalaire), induisant une **pénalisation quadratique**.
- $\sum_{j=1}^J$: cette régularisation est appliquée à toutes les feuilles $j = 1, \dots, J$.
- λ : hyperparamètre qui permet de réduire la complexité du modèle et d'éviter l'overfitting. λ contrôle la pénalisation des poids (plus λ est grand, plus on freine les valeurs extrêmes de γ_j).

- Pénalisation structurelle :

$$\gamma \sum_{k=1}^K d_k$$

Cette somme pénalise la complexité structurelle des arbres du modèle en limitant leur profondeur.

- Signification des symboles :

- d_k : profondeur ou complexité de l'arbre k , par exemple le nombre de feuilles.
- $\sum_{k=1}^K d_k$: somme de toutes les profondeurs des arbres $k = 1, \dots, K$.
- γ : hyperparamètre qui contrôle la pénalisation sur la taille des arbres (plus γ est élevé, plus le modèle favorise des arbres simples, moins profonds).

5.4.2 Optimisation de la fonction objectif via l'approximation de Taylor

L'une des contributions essentielles de XGBoost est l'utilisation d'une écriture de la fonction objectif permettant une optimisation efficace à chaque itération m . Cela repose sur un développement de Taylor à l'ordre 2 de la fonction de perte, qui facilite l'utilisation conjointe du gradient et de la hessienne pour accélérer la convergence.

Approximation de la fonction objectif

On cherche à minimiser la fonction objectif suivante à l'itération m :

$$\mathcal{L}^{(m)}(f_m) = \sum_{i=1}^N L(y_i, \hat{y}_i^{(m-1)} + f_m(x_i)) + \Omega(f_m)$$

En développant $L(y_i, \hat{y}_i^{(m-1)} + f_m(x_i))$ au voisinage de $\hat{y}_i^{(m-1)}$ avec un développement de Taylor à l'ordre 2, on obtient :

$$\mathcal{L}^{(m)}(f_m) \approx \sum_{i=1}^N \left[g_i \cdot f_m(x_i) + \frac{1}{2} h_i \cdot f_m^2(x_i) \right] + \Omega(f_m)$$

avec :

- $g_i = \left. \frac{\partial L(y_i, \hat{y}_i)}{\partial \hat{y}_i} \right|_{\hat{y}_i = \hat{y}_i^{(m-1)}}$: le gradient (dérivée première)
- $h_i = \left. \frac{\partial^2 L(y_i, \hat{y}_i)}{\partial \hat{y}_i^2} \right|_{\hat{y}_i = \hat{y}_i^{(m-1)}}$: la hessienne (dérivée seconde)

Structure de f_m et simplification de l'écriture

On suppose que $f_m(x_i)$ correspond au poids w_j de la feuille j de l'arbre dans laquelle tombe l'observation x_i . On note I_j l'ensemble des indices i tels que x_i est assigné à la feuille j , et J le nombre total de feuilles de l'arbre.

En remplaçant $f_m(x_i) = w_j$ pour $i \in I_j$, on obtient :

$$\mathcal{L}^{(m)} = \sum_{j=1}^J \left[\left(\sum_{i \in I_j} g_i \right) w_j + \frac{1}{2} \left(\sum_{i \in I_j} h_i + \lambda \right) w_j^2 \right] + \gamma J$$

Interprétation des termes de régularisation

- λ est un hyperparamètre qui pénalise les valeurs élevées des poids w_j , pour éviter des scores extrêmes et donc du surapprentissage.
- γ pénalise le nombre total de feuilles J , ce qui favorise les arbres plus simples.

Calcul du poids optimal w_j de chaque feuille

On minimise la fonction objectif par rapport à chaque w_j :

$$\frac{\partial \mathcal{L}^{(m)}}{\partial w_j} = \left(\sum_{i \in I_j} g_i \right) + \left(\sum_{i \in I_j} h_i + \lambda \right) w_j = 0$$

Ce qui donne la solution optimale :

$$w_j^* = - \frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}$$

Perte optimisée associée à un arbre

En remplaçant w_j^* dans l'expression de $\mathcal{L}^{(m)}$, on obtient la perte finale de l'arbre construit à l'itération m :

$$\mathcal{L}^{(m)} = -\frac{1}{2} \sum_{j=1}^J \frac{\left(\sum_{i \in I_j} g_i\right)^2}{\sum_{i \in I_j} h_i + \lambda} + \gamma J$$

Cette quantité est utilisée pour déterminer le gain d'information G lors du choix optimal de split dans l'arbre.

5.4.3 Critère de Split Optimal dans XGBoost

Pour qu'une séparation (split) d'une feuille en deux nouvelles feuilles soit considérée comme intéressante, il est nécessaire que la somme des pertes pénalisées des deux nouvelles feuilles soit inférieure à la perte de la feuille initiale, c'est-à-dire :

$$\text{LOSS}_{\text{Feuille Gauche}} + \text{LOSS}_{\text{Feuille Droite}} < \text{LOSS}_{\text{Feuille Initiale}}$$

L'expression de la perte totale pour un arbre construit à l'itération m , notée $\mathcal{L}^{(m)}$, est donnée par :

$$\mathcal{L}^{(m)} = -\frac{1}{2} \sum_{j=1}^J \frac{\left(\sum_{i \in I_j} g_i\right)^2}{\sum_{i \in I_j} h_i + \lambda} + \gamma J$$

Pour une feuille unique (avant split), la perte est donc :

$$-\frac{1}{2} \frac{\left(\sum_{i \in I_{\text{FI}}} g_i\right)^2}{\sum_{i \in I_{\text{FI}}} h_i + \lambda} + \gamma$$

où :

- I_{FI} représente les indices des observations dans la feuille initiale,
- I_{FG} et I_{FD} pour les feuilles gauche et droite.

Ainsi, le split est bénéfique si :

$$-\frac{1}{2} \frac{\left(\sum_{i \in I_{\text{FG}}} g_i\right)^2}{\sum_{i \in I_{\text{FG}}} h_i + \lambda} - \frac{1}{2} \frac{\left(\sum_{i \in I_{\text{FD}}} g_i\right)^2}{\sum_{i \in I_{\text{FD}}} h_i + \lambda} + 2\gamma < -\frac{1}{2} \frac{\left(\sum_{i \in I_{\text{FI}}} g_i\right)^2}{\sum_{i \in I_{\text{FI}}} h_i + \lambda} + \gamma$$

Ce qui peut se réécrire en simplifiant :

$$\frac{\left(\sum_{i \in I_{\text{FG}}} g_i\right)^2}{\sum_{i \in I_{\text{FG}}} h_i + \lambda} + \frac{\left(\sum_{i \in I_{\text{FD}}} g_i\right)^2}{\sum_{i \in I_{\text{FD}}} h_i + \lambda} - \frac{\left(\sum_{i \in I_{\text{FI}}} g_i\right)^2}{\sum_{i \in I_{\text{FI}}} h_i + \lambda} - 2\gamma > 0$$

Ainsi ce critère signifie qu'une séparation est avantageuse si elle réduit la perte globale, en tenant compte des pénalités de complexité. En pratique, XGBoost applique ce critère à chaque possible séparation d'une feuille, et choisit celle qui maximise ce gain de perte.

5.4.4 Parallélisation et Traitement des Valeurs Manquantes

XGBoost est conçu pour être rapide, même sur des ensembles de données volumineux, en tirant parti de la parallélisation. Les calculs de gradients et de hessiennes sont parallélisés, permettant une exécution plus rapide. De plus, XGBoost gère les **valeurs manquantes** de manière efficace, en décidant automatiquement de la meilleure façon de traiter ces valeurs lors de la construction des arbres.

5.5 Synthèse

En conclusion, XGBoost améliore le Gradient Boosting par l'introduction de la régularisation, de la descente de gradient second ordre, de la parallélisation, et d'optimisations sophistiquées dans la recherche des splits. Ces améliorations permettent à XGBoost de surpasser les méthodes traditionnelles en termes de rapidité et de performance, ce qui en fait un choix privilégié pour des tâches complexes de régression et de classification.

5.6 Application numérique dans le cadre d'une régression

Soit le jeu de données suivant, qui mesure l'efficacité d'un médicament en fonction de son dosage. Soit x le dosage et y la mesure de l'efficacité. L'objectif est de construire une suite d'arbres pour prédire en fonction de x la valeur de y :

X (dosage)	10	20	25	35
Y (efficacité)	-10	7	8	-7

Ce jeu de données est simpliste : 4 observations et 1 feature. Cela permettra d'illustrer parfaitement l'algorithme XGBoost pour la régression.

Étapes de l'algorithme

- Étape 0 : 1er estimateur constant $\hat{y}^{(0)}$
- Étape 1 : calcul des résidus $y_i - \hat{y}^{(0)}$
- Étape 2 : construction du premier arbre
 - identification du split optimal
 - calcul des estimations des feuilles
 - itération jusqu'à l'absence de split possible
- Étape 3 : calcul des résidus pour l'arbre suivant

Étape 0 : initialisation

- On prend $\hat{y}^{(0)} = 0.5$

Étape 1 : calcul des résidus

X (dosage)	10	20	25	35
Y (efficacité)	-10	7	8	-7
1 ^{er} est. $\hat{y}^{(0)}$	0.5	0.5	0.5	0.5
Résidus	-10.5	6.5	7.5	-7.5

Étape 2 : construction de l'arbre

Dans le cadre de la régression, nous choisissons d'utiliser la fonction de perte quadratique classique, définie comme suit :

$$L(y_i, \hat{y}_i) = \frac{1}{2} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

où \hat{y}_i est la prédiction du modèle pour l'observation i , et y_i est la valeur réelle.

Calcul du gradient

On souhaite dériver la fonction de perte par rapport à \hat{y}_i pour obtenir le gradient :

$$g_i = \frac{\partial}{\partial \hat{y}_i} L(y_i, \hat{y}_i) = \frac{\partial}{\partial \hat{y}_i} \left(\frac{1}{2} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \right) = (\hat{y}_i - y_i)$$

Ainsi, on obtient :

$$g_i = \hat{y}_i - y_i$$

Calcul du hessien

On dérive à nouveau pour obtenir la dérivée seconde :

$$h_i = \frac{\partial^2}{\partial \hat{y}_i^2} L(y_i, \hat{y}_i) = \frac{\partial}{\partial \hat{y}_i} (\hat{y}_i - y_i) = 1$$

Le hessien est donc constant :

$$h_i = 1 \quad (\text{valeur constante dans le cas quadratique})$$

Formule du score :

$$\text{Score} = \frac{(\sum g_i)^2}{\sum h_i + \lambda} + \frac{(\sum g_i)^2}{\sum h_i + \lambda} - \frac{(\sum g_i)^2}{\sum h_i + \lambda} - 2\gamma$$

Cas 1 : $\lambda = 0, \gamma = 0$.

Split du 1er niveau de l'arbre : comparons le score de chacun des splits. On a 4 valeurs dans cet exemple, donc 3 splits possibles.

	Split 1.1 (x < 15)	Split 1.2 (x < 22.5)	Split 1.3 (x < 30)
Feuille initiale – Résidus	[-10.5, 6.5, 7.5, -7.5]	[-10.5, 6.5, 7.5, -7.5]	[-10.5, 6.5, 7.5, -7.5]
Feuille initiale – Score ($\lambda = 0$)	$\frac{(-10.5 + 6.5 + 7.5 - 7.5)^2}{4} = 4$		
Feuille gauche – Résidus	[-10.5]	[-10.5, 6.5]	[-10.5, 6.5, 7.5]
Feuille gauche – Score ($\lambda = 0$)	$\frac{(-10.5)^2}{1} = 110.25$	$\frac{(-10.5 + 6.5)^2}{2} = 8$	$\frac{(-10.5 + 6.5 + 7.5)^2}{3} = 4.08$
Feuille droite – Résidus	[6.5, 7.5, -7.5]	[7.5, -7.5]	[-7.5]
Feuille droite – Score ($\lambda = 0$)	$\frac{(6.5 + 7.5 - 7.5)^2}{3} = 14.08$	$\frac{(7.5 - 7.5)^2}{2} = 0$	$\frac{(-7.5)^2}{1} = 56.25$
Total Score ($\lambda = 0, \gamma = 0$)	$110.25 + 14.08 - 4 = 120.33$	$8 + 0 - 4 = 4$	$4.08 + 56.25 - 4 = 56.33$

TABLE 1 – Comparaison des splits du premier niveau de l'arbre ($\lambda = 0, \gamma = 0$)

Le split 1.1 (x<15) a le plus fort score, il est donc conservé dans la construction de l'arbre au premier niveau

Maintenant nous allons spliter sur la feuille de droite.

On split la feuille de droite.

	Split 2.1 (x < 22.5)	Split 2.2 (x < 30)
Feuille initiale – Résidus	[6.5, 7.5, -7.5]	[6.5, 7.5, -7.5]
Feuille initiale – Score ($\lambda = 0$)	$\frac{(6.5 + 7.5 - 7.5)^2}{3} = 14.08$	
Feuille gauche - Résidus	[6.5]	[-6.5, 7.5]
Feuille gauche – Score ($\lambda = 0$)	$\frac{(6.5)^2}{1} = 42.25$	$\frac{(6.5 + 7.5)^2}{2} = 98$
Feuille droite - Résidus	[7.5, -7.5]	[-7.5]
Feuille droite – Score ($\lambda = 0$)	$\frac{(7.5 - 7.5)^2}{2} = 0$	$\frac{(-7.5)^2}{1} = 56.25$
Total Score ($\lambda = 0, \gamma = 0$)	$42.25 + 0 - 14.08 = 28.17$	$98 + 56.25 - 14.08 = 140.17$

TABLE 2 – Scores des splits du deuxième niveau de l'arbre (cas $\lambda = 0, \gamma = 0$)

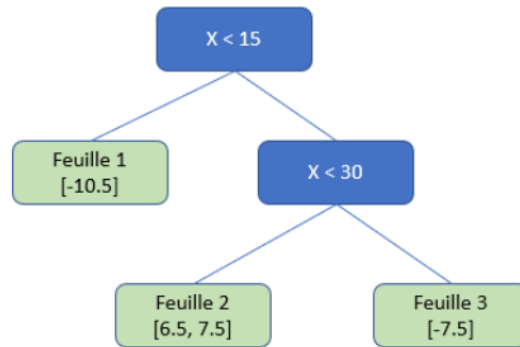


FIGURE 10 – Arbre obtenu

Le split 2.2 ($x < 30$) a le plus fort gain, il est donc conservé dans la construction de l'arbre. On suppose maintenant que notre arbre est complètement construit, on calcule maintenant les estimations optimales par feuille. On se sert des résidus restant par feuille.

Calcul des sorties par feuille

$$w_j = \frac{-\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}$$

Feuille	Résidus	Output ($\lambda = 0$)
Feuille 1	[-10.5]	-10.5
Feuille 2	[6.5, 7.5]	7
Feuille 3	[-7.5]	-7.5

Nouvelle estimation :

$$\hat{y}^{(1)} = \hat{y}^{(0)} + f_1(x_i)$$

X	10	20	25	35
Y	-10	7	8	-7
$\hat{y}^{(0)}$	0.5	0.5	0.5	0.5
Résidu	-10.5	6.5	7.5	-7.5
$f_1(x_i)$	-10.5	7	7	-7.5
$\hat{y}^{(1)}$	-10	7.5	7.5	-7

Cas 2 : $\lambda = 1, \gamma = 0$

Split du 1er niveau de l'arbre :

Split du 1 ^{er} niveau de l'arbre	Split 1.1 ($x < 15$)	Split 1.2 ($x < 22.5$)	Split 1.3 ($x < 30$)
Feuille initiale – Résidus	$[-10.5, 6.5, 7.5, -7.5]$	$[-10.5, 6.5, 7.5, -7.5]$	$[-10.5, 6.5, 7.5, -7.5]$
Feuille initiale – Score ($\lambda = 1$)	$\frac{(-10.5 + 6.5 + 7.5 - 7.5)^2}{4 + 1} = 3.2$		
Feuille gauche – Résidus	$[-10.5]$	$[-10.5, 6.5]$	$[-10.5, 6.5, 7.5]$
Feuille gauche – Score ($\lambda = 1$)	$\frac{(-10.5)^2}{2} = 55.12$	$\frac{(-10.5 + 6.5)^2}{3} = 5.33$	$\frac{(-10.5 + 6.5 + 7.5)^2}{4} = 3.06$
Feuille droite – Résidus	$[6.5, 7.5, -7.5]$	$[7.5, -7.5]$	$[-7.5]$
Feuille droite – Score ($\lambda = 1$)	$\frac{(6.5 + 7.5 - 7.5)^2}{4} = 10.56$	$\frac{(7.5 - 7.5)^2}{3} = 0$	$\frac{(-7.5)^2}{2} = 28.1$
Total Score ($\lambda = 1, \gamma = 0$)	$55.12 + 10.56 - 3.2 = 62.48$	$5.33 - 3.2 = 2.13$	$3.06 + 28.1 - 3.2 = 27.96$

TABLE 3 – Comparaison des splits au 1^{er} niveau de l'arbre avec pénalisation $\lambda = 1$

On voit que c'est toujours le split 1.1 ($x < 15$) qui a le plus fort score, il est donc conservé dans la construction de l'arbre au premier niveau. Par contre, l'ajout d'une pénalisation $\lambda = 1$ abaisse fortement les scores, ce qui amène à des arbres beaucoup moins profonds que sans pénalisation. C'est bien le but recherché : avoir des arbres moins profonds permettant de limiter les risques de sur apprentissage.

Nous obtenons les mêmes splits que pour le premier cas. Maintenant en calculant les sorties par feuille on obtient :

Feuille	Résidus	Output ($\lambda = 0$)	Output ($\lambda = 1$)
Feuille 1	$[-10.5]$	$\frac{-10.5}{1} = -10.5$	$\frac{-10.5}{2} = -5.25$
Feuille 2	$[6.5, 7.5]$	$\frac{6.5 + 7.5}{2} = 7$	$\frac{6.5 + 7.5}{3} = 4.6$
Feuille 3	$[-7.5]$	$\frac{-7.5}{1} = -7.5$	$\frac{-7.5}{2} = -3.75$

TABLE 4 – Estimations optimales par feuille pour $\lambda = 0$ et $\lambda = 1$

Nouvelle estimation :

$$\hat{y}^{(1)} = \hat{y}^{(0)} + f_1(x_i)$$

X	10	20	25	35
Y	-10	7	8	-7
$\hat{y}^{(0)}$	0.5	0.5	0.5	0.5
Résidu	-10.5	6.5	7.5	-7.5
$f_1(x_i)$	-5.25	4.6	4.6	3.75
$\hat{y}^{(1)}$	-4.75	5.1	5.1	-3.25

Cas 3 : $\gamma > 0$

$$\text{Score} = \text{Score}_{FG} + \text{Score}_{FD} - \text{Score}_{FI} - 2\gamma$$

Avec $\gamma > 0$ cela réduit les possibilités qu'un split soit valide (il faut que le score soit positif). Et donc par conséquent, cela limite la profondeur des arbres, et donc les risques de sur-apprentissage.

6 Simulation

6.1 Cadre de la simulation

Nous allons maintenant présenter la simulation que nous avons mise en place afin de pouvoir évaluer les performances du modèle CART et des algorithmes **AdaBoost** et **XGBoost** dans le cadre de la détection des maladies cardiaques grâce à différentes variables.

Notre simulation repose sur un jeu de données regroupant plus de 900 cas comportant une maladie ou non. Pour cela nous avons mis en place une classification supervisée sur la présence ou non de maladie cardiaque. Pour ce qui est de nos variables explicatives nous avons :

- des variables biologiques avec l'âge du patient et son sexe (encodé : 1=homme et 0=femme)
- des variables cliniques avec la pression artérielle(en mm Hg) au repos, le taux de cholestérol ((en mg/dL), la glycémie à jeun (encodé 1 si > 120 mg/dL, sinon 0), la fréquence cardiaque maximale pendant un effort et la dépression du segment ST induite par l'exercice (Le segment ST est une partie de l'électrocardiogramme qui représente la période entre la dépolarisation et la repolarisation des ventricules du cœur. Ainsi une dépression de ce segment le muscle cardiaque ne reçoit pas assez d'oxygène à ce moment-là)
- une variable de douleur à l'effort avec la présence ou non d'angine de poitrine provoquée par l'effort.
- des variables analysant les douleurs thoraciques (angineuse ou non)
- une variable d'analyse des électrocardiogrammes des patients
- des variables analysant la pente du segment ST (plate ou montante)

Nous avons ensuite entraîné nos trois modèles de prédiction sur ces variables à travers un code Python (cf Annexe) et affiché les métriques de prédiction de nos modèles ainsi que les matrices de confusion.

Rappel sur les métriques d'évaluation

Pour évaluer les performances de nos modèles de classification, nous utilisons les métriques suivantes :

- **Précision** : proportion de vraies prédictions positives parmi toutes les prédictions positives.

$$\text{Précision} = \frac{VP}{VP + FP}$$

où VP est le nombre de vrais positifs et FP le nombre de faux positifs.

- **Rappel** : proportion de vraies prédictions positives parmi toutes les instances positives réelles.

$$\text{Rappel} = \frac{VP}{VP + FN}$$

où FN est le nombre de faux négatifs.

- **F1-score** : moyenne harmonique entre la précision et le rappel, utile en cas de classes dés-équilibrées.

$$F1 = 2 \times \frac{\text{Précision} \times \text{Rappel}}{\text{Précision} + \text{Rappel}}$$

6.2 Comportement modèle CART

Le modèle a été testé sur 368 données (40% de test et 60% d'entraînement) et voici ces métriques de performance par classe :

Classe	Précision	Rappel	F1-score	Nombre d'exemples
0 (pas de maladie)	0.76	0.75	0.76	173
1 (maladie cardiaque)	0.78	0.79	0.79	195

TABLE 5 – Métriques de performance du modèle CART par classe

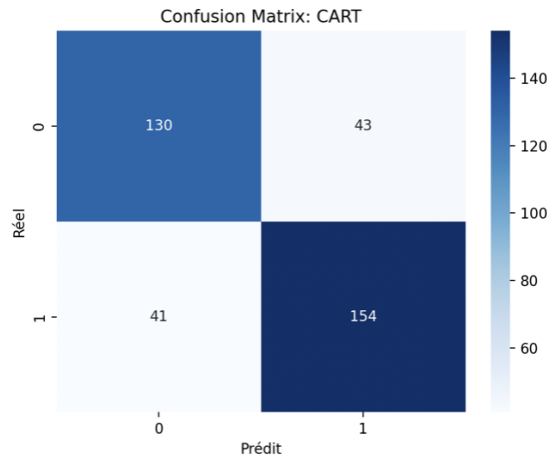


FIGURE 11 – Matrice de confusion du modèle CART

Le modèle parvient à équilibrer correctement les prédictions entre les deux classes. Le nombre d'erreurs de prédiction reste relativement similaire pour chaque classe : on observe 41 faux négatifs et 43 faux positifs. Cela montre que le modèle ne favorise pas une classe au détriment de l'autre. Le F1-score moyen atteint 0,77, ce qui reflète une performance globale raisonnable et équilibrée.

6.3 Comportement d'AdaBoost

Comme pour le modèle **CART**, **AdaBoost** a été testé sur 368 données (40% de test et 60% d'entraînement) et voici ces métriques de performance par classe :

Classe (réelle)	Précision	Rappel	F1-score	Nombre d'exemples
0 (pas de maladie)	0.82	0.83	0.82	173
1 (maladie cardiaque)	0.85	0.84	0.84	195

TABLE 6 – Métriques de performance du modèle AdaBoost par classe

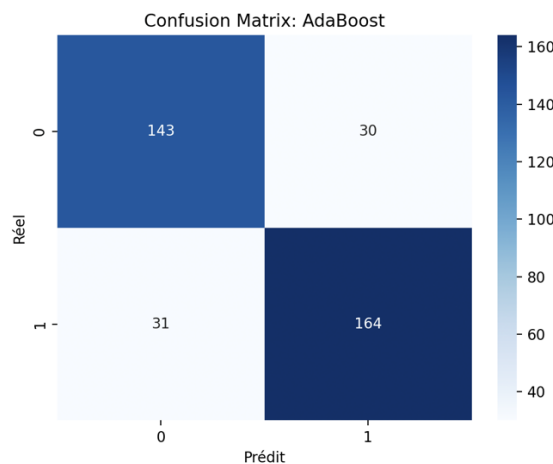


FIGURE 12 – Matrice de confusion d'AdaBoost

Bien que l'on observe un léger déséquilibre entre les prédictions des deux classes, AdaBoost montre de meilleures performances que CART. Il obtient une précision de 0,82 pour la classe 0 et de 0,85 pour la classe 1. La matrice de confusion indique une réduction du nombre d'erreurs : 30 faux positifs et 31 faux négatifs. Globalement, le F1-score moyen de 0,83 témoigne d'une amélioration par rapport au modèle CART, notamment en termes de précision et de rappel pour chaque classe.

6.4 Comportement de XGBoost

Comme pour les deux précédents cas, **XGBoost** a été testé sur 368 données (40% de test et 60% d'entraînement) et voici ces métriques de performance par classe :

Classe (réelle)	Précision	Rappel	F1-score	Nombre d'exemples
0 (pas de maladie)	0.86	0.82	0.84	173
1 (maladie cardiaque)	0.85	0.88	0.86	195

TABLE 7 – Métriques de performance du modèle XGBoost par classe

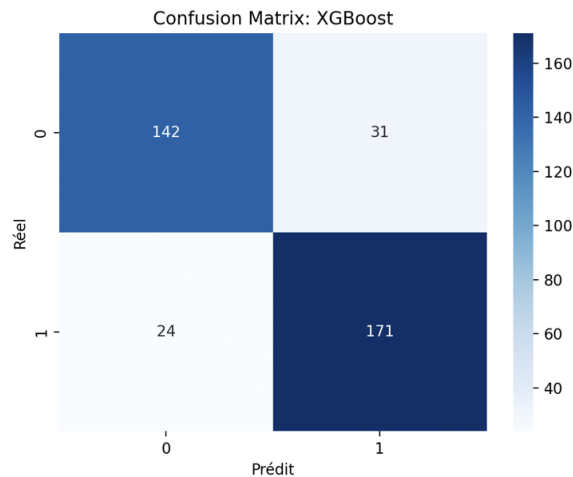


FIGURE 13 – Matrice de confusion pour XGBoost

Nous remarquons que le modèle XGBoost affiche des performances légèrement supérieures à **AdaBoost**, notamment avec un F1-score global de 0,85. Ce bon résultat s'explique par un équilibre efficace entre la précision et le rappel pour chaque classe.

Pour la classe 0 (absence de maladie), la précision est de 0,86 et le rappel de 0,82, ce qui traduit une bonne capacité du modèle à éviter les faux positifs. Pour la classe 1 (présence d'une maladie cardiaque), le rappel est élevé à 0,88, ce qui signifie que XGBoost détecte efficacement les cas de patients malades, ce qui est crucial dans un contexte médical.

La matrice de confusion montre que le modèle commet 24 faux négatifs et 31 faux positifs, ce qui est meilleur que le modèle AdaBoost, notamment sur les faux négatifs, ce qui diminue le risque de ne pas diagnostiquer une maladie existante.

6.5 Comparaison globale

Après avoir analysé séparément les performances de chaque modèle (**CART**, **AdaBoost** et **XGBoost**), voyons ici une synthèse comparative. Les principales métriques de classification (précision, rappel, F1-score et exactitude) sont rassemblées dans le tableau suivant :

Modèle	Précision	Rappel	F1-Score	Accuracy
CART	0.77	0.77	0.77	0.77
AdaBoost	0.83	0.83	0.83	0.83
XGBoost	0.85	0.85	0.85	0.85

TABLE 8 – Résultats des modèles sur la prédiction de la maladie cardiaque

L'évaluation conjointe des trois modèles, CART, AdaBoost et XGBoost, met en évidence des différences notables en termes de performances, de stabilité et de pertinence dans un contexte de prédiction médicale.

Le modèle CART, basé sur des arbres de décision simples, constitue une solution rapide et interprétable. Toutefois, ses performances restent limitées avec un F1-score de 0,77. Il tend à générer davantage de faux positifs et de faux négatifs, ce qui peut poser problème dans une tâche sensible comme la détection de maladies.

L'introduction du boosting avec AdaBoost permet une nette amélioration. En combinant plusieurs arbres faibles, le modèle parvient à corriger progressivement les erreurs du modèle de base. Cette approche se traduit par une hausse significative des performances (F1-score de 0,83) et une réduction des erreurs critiques. AdaBoost montre ainsi sa capacité à mieux généraliser sans surapprentissage excessif.

Enfin, XGBoost pousse cette logique encore plus loin, grâce à une régularisation fine et une optimisation efficace de la fonction objective. Il affiche les meilleures performances (F1-score de 0,85) tout en maintenant un bon équilibre entre précision et rappel. De plus, il parvient à réduire davantage les faux négatifs, ce qui est essentiel dans le cadre d'un diagnostic médical où les erreurs d'omission peuvent avoir des conséquences graves.

En croisant les résultats, on observe une progression cohérente : chaque modèle apporte un gain par rapport au précédent, aussi bien en précision qu'en robustesse. CART pose les bases, AdaBoost les améliore, et XGBoost les consolide. Ce constat justifie pleinement le recours aux modèles d'ensemble pour les problèmes de classification médicale, en particulier lorsque la sensibilité du modèle (rappel) est un enjeu majeur.

7 Conclusion

Ce travail de recherche a permis d'explorer en profondeur les fondements, les mécanismes et les performances des algorithmes de boosting appliqués à l'apprentissage supervisé, que ce soit dans un contexte de régression ou de classification. En partant du modèle de base CART, nous avons mis en évidence les limites des arbres de décision lorsqu'ils sont utilisés isolément, notamment leur sensibilité au surapprentissage et leur manque de robustesse sur des jeux de données complexes ou bruités.

L'introduction d'AdaBoost a marqué une avancée significative. En combinant plusieurs classificateurs faibles de manière séquentielle et pondérée, cet algorithme améliore considérablement la capacité de généralisation du modèle. Nous avons démontré que son fondement repose sur une interprétation probabiliste et une minimisation de la perte exponentielle. Toutefois, nous avons également mis en lumière certaines de ses limites, notamment sa sensibilité aux données aberrantes.

Pour pallier ces limites, XGBoost s'impose comme une solution de pointe. En capitalisant sur les principes du gradient boosting, tout en y ajoutant des techniques avancées telles que la régularisation, l'optimisation via un développement de Taylor et une gestion intelligente des données manquantes, XGBoost se démarque à la fois par sa performance statistique et son efficacité computationnelle. Les expérimentations menées dans le cadre de ce travail confirment la supériorité de ce modèle, notamment dans des contextes où la minimisation des erreurs critiques comme les faux négatifs est primordiale, à l'instar des applications médicales.

Il est important de souligner que ce travail ne s'est pas limité à une étude empirique des algorithmes. Il s'est également appuyé sur une formalisation théorique rigoureuse du boosting, en s'appuyant notamment sur le théorème de Schapire et la modélisation additive progressive. Ces fondements ont permis de mieux comprendre pourquoi les algorithmes comme AdaBoost ou XGBoost parviennent à améliorer la performance des classificateurs faibles, en concentrant progressivement l'apprentissage sur les erreurs commises par les modèles précédents. Cette articulation entre théorie et pratique a constitué l'un des fils conducteurs du mémoire.

Plus globalement, ce TER a été l'occasion de mettre en pratique les concepts théoriques de l'apprentissage statistique à travers une étude comparative rigoureuse et des implémentations concrètes. Il illustre de manière concrète comment les avancées algorithmiques récentes, notamment dans le domaine du boosting, permettent de relever efficacement les défis posés par la classification supervisée.

Annexe

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Nous d finissons un vecteur de marges yf(x)
5 marge = np.linspace(-2, 2, 400)
6
7 # Nous d finissons les fonctions de perte
8 def perte_exponentielle(yf):
9     return np.exp(-yf)
10
11 def perte_deviance(yf):
12     return np.log(1 + np.exp(-2 * yf))
13
14 # Calcul des pertes
15 loss_exp = perte_exponentielle(marge)
16 loss_dev = perte_deviance(marge)
17
18 # Trac des graphiques
19 plt.figure(figsize=(8, 5))
20 plt.plot(marge, loss_exp, label="Exponentielle", color="skyblue", linewidth=2)
21 plt.plot(marge, loss_dev, label="D viance binomiale", color="orange", linewidth=2)
22
23 # Mise en forme
24 plt.xlabel("Marge $yf(x)$")
25 plt.ylabel("Perte")
26 plt.title("Comparaison des fonctions de perte")
27 plt.legend()
28 plt.tight_layout()
29 plt.show()
```

Listing 1 – Code Python utilisé pour tracer les fonctions de perte exponentielle et déviance binomiale en fonction de la marge

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import classification_report, confusion_matrix
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.ensemble import AdaBoostClassifier
6 from xgboost import XGBClassifier
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9
10 # Chargement du jeu de donn es
11 df = pd.read_csv("heart.csv")
12
13 print(df.columns.tolist())
14
15 # S lection manuelle des variables explicatives
16 features = [
17     'Age',
18     'RestingBP', # Pression art rielle au repos (mm/Hg)
19     'Cholesterol', # Taux de cholest rol s rique (mg/dL)
20     'FastingBS', # Glyc mie jeun
21     'MaxHR', # Fr quence cardiaque maximale l'effort
22     'Oldpeak', # D pression du segment ST (exercice)
23     'Sex_M',
24     'ChestPainType_ATA', #Douleur thoracique type ATA
25     'ChestPainType_NAP', #Douleur thoracique type NAP
26     'ChestPainType_TA', #Douleur thoracique type TA
27     'RestingECG_Normal', #Electrocardiogramme normal
28     'RestingECG_ST', #Electrocardiogramme avec anomalie
29     'ExerciseAngina_Y', #Douleur angineuse pendant l'effort
30     'ST_Slope_Flat', #Pente du segment ST plate
31     'ST_Slope_Up', #Pente du segment ST ascendante
32 ]
33
34 # X = variables explicatives, Y = cible (maladie cardiaque)
35 X = df[features]
```

```

36 Y = df['HeartDisease']
37
38 # S paration des donn es (60% entra nement, 40% test)
39 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.4, random_state
    =14)
40
41 # Fonction d' valuation
42 def eval_model(Y_true, Y_pred, title):
43     print(f"--- {title} ---")
44     print(classification_report(Y_true, Y_pred))
45     sns.heatmap(confusion_matrix(Y_true, Y_pred), annot=True, fmt='d', cmap='Blues')
46     plt.title(f"Confusion Matrix: {title}")
47     plt.xlabel("Pr dit")
48     plt.ylabel("R el")
49     plt.show()
50
51 # Mod le CART
52 cart = DecisionTreeClassifier(random_state=14)
53 cart.fit(X_train, Y_train)
54 Y_pred_cart = cart.predict(X_test)
55 eval_model(Y_test, Y_pred_cart, "CART")
56
57 # Mod le AdaBoost
58 ada = AdaBoostClassifier(n_estimators=100, random_state=14)
59 ada.fit(X_train, Y_train)
60 Y_pred_ada = ada.predict(X_test)
61 eval_model(Y_test, Y_pred_ada, "AdaBoost")
62
63 # Mod le XGBoost
64 xgb = XGBClassifier(use_label_encoder=False, eval_metric='logloss', random_state=14)
65 xgb.fit(X_train, Y_train)
66 Y_pred_xgb = xgb.predict(X_test)
67 eval_model(Y_test, Y_pred_xgb, "XGBoost")

```

Listing 2 – Code Python pour entraîner et évaluer les modèles CART, AdaBoost et XGBoost sur les données cardiaques

Références

- [1] Trevor Hastie, Robert Tibshirani, Jerome Friedman. *The Elements of Statistical Learning : Data Mining, Inference, and Prediction*. Second Edition, Springer, 2009.
- [2] StatQuest with Josh Starmer. *Regression trees clearly explained*. Youtube 2019. Disponible à l'adresse : https://www.youtube.com/watch?v=g9c66TUy1Z4&t=41s&ab_channel=StatQuestwithJoshStarmer
- [3] StatQuest with Josh Starmer. *Decision and classification trees clearly explained*. Youtube 2021. Disponible à l'adresse : https://www.youtube.com/watch?v=_L39rN6gz7Y&ab_channel=StatQuestwithJoshStarmer
- [4] Robert E. Schapire. *The strength of weak learnability*. Volume 5, pages 197-227 Springer, 1990. Disponible à l'adresse : <https://link.springer.com/article/10.1007/BF00116037>
- [5] Ricco Rakotomalala. *Gradient Boosting – Technique ensembliste pour l'analyse prédictive*. Université Lumière Lyon 2, 2020. Disponible à l'adresse : https://eric.univ-lyon2.fr/ricco/cours/slides/gradient_boosting.pdf
- [6] Objectif Data Science avec Vincent. *XGBoost : l'algorithme star de machine learning décortiqué*. YouTube, 2022. Disponible à l'adresse : https://www.youtube.com/watch?v=iD_TL725wZs
- [7] StatQuest with Josh Starmer. *XGBoost Part 1 (of 4) : Régression*. YouTube, 2020. Disponible à l'adresse : <https://www.youtube.com/watch?v=0tD8wVaFm6E>
- [8] StatQuest with Josh Starmer. *XGBoost Part 2 : Arbres XGBoost pour la classification*. YouTube, 2020. Disponible à l'adresse : <https://www.youtube.com/watch?v=8b1JEDvenQU>
- [9] Kaggle. *Heart Failure Prediction Dataset*, consulté en mai 2025. Disponible à l'adresse : <https://www.kaggle.com/code/tanmay111999/heart-failure-prediction-cv-score-90-5-models/input>