

**Sélection de modèles et sélection d'estimateurs pour
l'Apprentissage statistique (Cours Peccot)
Troisième cours: Rééchantillonnage et pénalisation**

SYLVAIN ARLOT
CNRS – ÉQUIPE SIERRA

TABLE DES MATIÈRES

1. Régressogrammes en régression hétéroscédastique	2
1.1. Cadre de la régression hétéroscédastique	2
1.2. Régressogrammes	2
1.3. Exemples de collections de partitions	4
1.4. Pénalité idéale	4
2. Nécessité d'estimer la forme de la pénalité	7
2.1. Illustration sur un exemple	7
2.2. Caractérisation des pénalités fonction de la dimension	8
2.3. Sous-optimalité des pénalités fonction de la dimension	9
2.4. Pourquoi estimer la forme de la pénalité?	9
3. Rééchantillonnage	10
3.1. Heuristique d'Efron	10
3.2. Rééchantillonnage à poids échangeables	11
3.3. Arguments théoriques asymptotiques	12
3.4. Usages du rééchantillonnage	13
3.5. Limites du rééchantillonnage	13
4. Un estimateur de la variance par rééchantillonnage	14
4.1. Formule close	14
4.2. Comparaison avec l'estimateur classique	15
4.3. Rééchantillonnage et structure	15
4.4. Rééchantillonnage et concentration	17
4.5. Calcul de la constante multiplicative	18
5. Pénalités par rééchantillonnage	19
6. Garanties théoriques pour les régressogrammes	20
6.1. Calcul de la pénalité par rééchantillonnage	21
6.2. Inégalité-oracle	24
6.3. Adaptation	25
6.4. Comparaison des poids	28
7. Estimation de densité par moindres carrés	28

1. RÉGRESSOGRAMMES EN RÉGRESSION HÉTÉROSCÉDASTIQUE

Une bonne partie des résultats de ce cours et du suivant portent sur le problème de sélection parmi des régressogrammes en régression hétéroscédastique sur un plan d'expérience aléatoire. Commençons par décrire formellement ce cadre.

1.1. Cadre de la régression hétéroscédastique. On considère le cadre de la régression sur un plan d'expérience aléatoire, où l'on observe $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \mathbb{R}$ i.i.d. de loi commune P . Si $(X, Y) \sim P$, on note

$$Y = \eta(X) + \varepsilon$$

où la fonction de régression $\eta : \mathcal{X} \mapsto \mathbb{R}$ est définie par $\eta(X) = \mathbb{E}[Y | X]$ P -p.s. et le bruit ε vérifie donc P -p.s.

$$\mathbb{E}[\varepsilon | X] = 0 \quad \text{et} \quad \mathbb{E}[\varepsilon^2 | X] = \sigma^2(X) .$$

La fonction inconnue $\sigma : \mathcal{X} \mapsto [0, +\infty[$ mesure le niveau de bruit des observations en fonction de leur position. On parle de *régression hétéroscédastique* dans la mesure où la distribution conditionnelle de ε sachant X est autorisée à varier en fonction de X , en particulier via le niveau de bruit $\sigma(\cdot)$.

Notons ici l'intérêt fondamental de supposer le plan d'expérience aléatoire : cela nous permet de conserver un échantillon $D_n = (X_i, Y_i)_{1 \leq i \leq n}$ de données i.i.d., alors que si le plan d'expérience était déterministe, les résidus $\varepsilon_1, \dots, \varepsilon_n$ ne seraient pas identiquement distribués.

Rappelons enfin que l'objectif est de construire un prédicteur $t \in \mathbb{S}$ (l'ensemble des applications mesurables $\mathcal{X} \mapsto \mathbb{R}$) dont la perte quadratique

$$P\gamma(t) = \mathbb{E}_{(X,Y) \sim P} [\gamma(t; (X, Y))] = \mathbb{E}_{(X,Y) \sim P} \left[(t(X) - Y)^2 \right]$$

est minimale, avec $\gamma : \mathcal{X} \times \mathbb{R} \mapsto \mathbb{R}$ le contraste des moindres carrés défini par

$$\forall (x, y) \in \mathcal{X} \times \mathbb{R} , \quad \forall t \in \mathbb{S} , \quad \gamma(t; (x, y)) = (t(x) - y)^2 .$$

La perte est minimale pour $t = s^*$ et l'on définit alors la perte relative

$$\ell(s^*, t) = P\gamma(t) - P\gamma(s^*) = \mathbb{E}_{(X,Y) \sim P} \left[(s^*(X) - t(X))^2 \right] .$$

1.2. Régressogrammes. Pour toute partition finie m de \mathcal{X} , on définit l'ensemble des fonctions constantes sur chaque élément de m (ou «modèle des histogrammes associé à m ») par

$$S_m := \left\{ t = \sum_{\lambda \in m} \alpha_\lambda \mathbb{1}_\lambda \text{ t.q. } \alpha \in \mathbb{R}^m \right\} .$$

Un estimateur des moindres carrés sur S_m (ou *régressogramme associé à la partition m*) est alors défini par

$$\widehat{s}_m = \widehat{s}_m(D_n) \in \operatorname{argmin}_{t \in S_m} \{P_n \gamma(t)\} \quad \text{où} \quad P_n = \frac{1}{n} \sum_{i=1}^n \delta_{(X_i, Y_i)}$$

est la mesure empirique associée à l'échantillon D_n .

Pour tout échantillon $D_n = (X_i, Y_i)_{1 \leq i \leq n}$ et tout $\lambda \subset \mathcal{X}$, on pose

$$\widehat{p}_\lambda = \widehat{p}_\lambda(D_n) = \frac{1}{n} \operatorname{Card} \{i \text{ t.q. } X_i \in \lambda\} .$$

Le régressogramme \widehat{s}_m associé à m est défini de manière unique si et seulement si

$$\min_{\lambda \in m} \{\widehat{p}_\lambda\} > 0 ,$$

et dans ce cas

$$\widehat{s}_m = \sum_{\lambda \in m} \widehat{\beta}_\lambda \mathbb{1}_\lambda \quad \text{avec} \quad \forall \lambda \in m , \quad \widehat{\beta}_\lambda = \widehat{\beta}_\lambda(D_n) := \frac{1}{n \widehat{p}_\lambda} \sum_{i \text{ t.q. } X_i \in \lambda} Y_i . \quad (1)$$

Preuve de l'Équation (1). Pour tout $t = \sum_{\lambda \in m} t_\lambda \mathbb{1}_\lambda \in S_m$, on a

$$P_n \gamma(t) = \frac{1}{n} \sum_{\lambda \in m} \left[\sum_{X_i \in \lambda} (Y_i - t_\lambda)^2 \right] . \quad (2)$$

Pour tout $\lambda \in m$, si $\widehat{p}_\lambda > 0$, $t_\lambda \mapsto \sum_{X_i \in \lambda} (Y_i - t_\lambda)^2$ admet un unique minimum sur \mathbb{R} en $t_\lambda = \widehat{\beta}_\lambda$, ce qui prouve (1) lorsque $\min_{\lambda \in m} \{\widehat{p}_\lambda\} > 0$. À l'inverse, s'il existe un $\lambda \in m$ pour lequel $\widehat{p}_\lambda = 0$, le risque empirique de t ne dépend pas de t_λ , et donc \widehat{s}_m n'est pas défini de manière unique sur chaque $\lambda \in m$ pour lequel $\widehat{p}_\lambda = 0$. \square

De la même manière, pour tout $t = \sum_{\lambda \in m} t_\lambda \mathbb{1}_\lambda \in S_m$, on a

$$\begin{aligned} \ell(s^*, t) &= \mathbb{E} \left[(t(X) - s^*(X))^2 \right] \\ &= \sum_{\lambda \in m} \left(p_\lambda \mathbb{E} \left[(t_\lambda - s^*(X))^2 \mid X \in \lambda \right] \right) \end{aligned} \quad (3)$$

ce qui prouve que

$$s_m^*(P) = s_m^* \in \operatorname{argmin}_{t \in S_m} \{P \gamma(t)\}$$

existe toujours, et est unique si et seulement si

$$\min_{\lambda \in m} p_\lambda > 0 \quad \text{où} \quad \forall \lambda \in m , \quad p_\lambda := \mathbb{P}_{(X,Y) \sim P} (X \in \lambda) .$$

Dans ce cas, on a

$$s_m^* = \sum_{\lambda \in m} \beta_\lambda \mathbb{1}_\lambda \quad \text{avec} \quad \forall \lambda \in m , \quad \beta_\lambda := \mathbb{E}_{(X,Y) \sim P} [Y \mid X \in \lambda] \quad (4)$$

et l'erreur d'approximation du modèle S_m vaut

$$\ell(s^*, s_m^*) = \sum_{\lambda \in m} \left(p_\lambda \left(\sigma_\lambda^{(d)} \right)^2 \right) \quad (5)$$

$$\text{où } \left(\sigma_\lambda^{(d)} \right)^2 := \mathbb{E} \left[(\beta_\lambda - s^*(X))^2 \mid X \in \lambda \right] .$$

1.3. Exemples de collections de partitions. On peut notamment considérer les deux exemples suivants de partitions de $\mathcal{X} = [0, 1[$:

- (1) Collection des partitions régulières $\mathcal{M}^{(\text{reg})} = \mathcal{M}^{(\text{reg})}(M_n)$, avec un nombre de morceaux maximal $M_n \geq 1$:

$$\mathcal{M}^{(\text{reg})} := \{ m_D \text{ t.q. } 1 \leq D \leq M_n \} \quad \text{avec} \quad m_D := \left(\left[\frac{k-1}{D}, \frac{k}{D} \right] \right)_{1 \leq k \leq D} .$$

- (2) Collection des partitions bi-régulières avec rupture en $1/2$ $\mathcal{M}_n^{(\text{reg}, 1/2)}$, avec un nombre de morceaux maximal $M_n \geq 2$:

$$\mathcal{M}^{(\text{reg})} := \{ m_1 \} \cup \left\{ m_{(D_1, D_2)} \text{ t.q. } 1 \leq D_1, D_2 \leq \frac{M_n}{2} \right\}$$

$$\text{avec} \quad m_{(D_1, D_2)} := \left(\left[\frac{k-1}{2D_1}, \frac{k}{2D_1} \right] \right)_{1 \leq k \leq D_1} \cup \left(\left[\frac{1}{2} + \frac{k-1}{2D_2}, \frac{1}{2} + \frac{k}{2D_2} \right] \right)_{1 \leq k \leq D_2} .$$

Lorsque $\mathcal{X} = [0, 1[$ ^{d} , on peut généraliser la collection $\mathcal{M}^{(\text{reg})}(M_n)$ comme suit :

$$\mathcal{M}^{(\text{reg})} := \left\{ m_T \text{ t.q. } 1 \leq T \leq M_n^{1/d} \right\} \quad \text{avec} \quad m_T := \left(\prod_{i=1}^d \left[\frac{k_i-1}{T}, \frac{k_i}{T} \right] \right)_{1 \leq k_1, \dots, k_d \leq T} .$$

1.4. Pénalité idéale. Nous sommes désormais en mesure de calculer la perte relative et la pénalité idéale associées au modèle S_m . En utilisant les notations de la Section 3.3 du deuxième cours, il s'agit de calculer les termes p_1 et p_2 .

On suppose désormais que $\min_{\lambda \in m} p_\lambda > 0$ (sinon, il suffit de fusionner les $\lambda \in m$ tels que $p_\lambda = 0$ avec un autre élément de m , sans rien changer au régressogramme correspondant).

Formules closes. Tout d'abord, d'après (3) et (4), pour tout $t = \sum_{\lambda \in m} t_\lambda \mathbb{1}_\lambda \in S_m$,

$$\begin{aligned} P(\gamma(t) - \gamma(s_m^*)) &= \sum_{\lambda \in m} \left[p_\lambda \left(\mathbb{E} \left[(t_\lambda - s^*(X))^2 \mid X \in \lambda \right] - \mathbb{E} \left[(\beta_\lambda - s^*(X))^2 \mid X \in \lambda \right] \right) \right] \\ &= \sum_{\lambda \in m} \left[p_\lambda (t_\lambda - \beta_\lambda)^2 \right] . \end{aligned}$$

Ainsi, en supposant que $\min_{\lambda \in m} \{\hat{p}_\lambda\} > 0$,

$$p_1(m) = P(\gamma(\hat{s}_m) - \gamma(s_m^*)) = \sum_{\lambda \in m} \left[p_\lambda \left(\hat{\beta}_\lambda - \beta_\lambda \right)^2 \right] . \quad (6)$$

De la même manière, d'après (2) et (1), pour tout $t = \sum_{\lambda \in m} t_\lambda \mathbb{1}_\lambda \in S_m$,

$$\begin{aligned} P_n(\gamma(t) - \gamma(\widehat{s}_m)) &= \frac{1}{n} \sum_{\lambda \in m} \sum_{X_i \in \lambda} \left[(t_\lambda - Y_i)^2 - (\widehat{\beta}_\lambda - Y_i)^2 \right] \\ &= \frac{1}{n} \sum_{\lambda \in m} \sum_{X_i \in \lambda} (t_\lambda - \widehat{\beta}_\lambda)^2 = \sum_{\lambda \in m} \widehat{p}_\lambda (t_\lambda - \widehat{\beta}_\lambda)^2 \end{aligned}$$

si bien que

$$p_2(m) = P_n(\gamma(s_m^*) - \gamma(\widehat{s}_m)) = \sum_{\lambda \in m} \left[\widehat{p}_\lambda (\widehat{\beta}_\lambda - \beta_\lambda)^2 \right]. \quad (7)$$

Notons que (7) est encore valable lorsque $\min_{\lambda \in m} \widehat{p}_\lambda = 0$ en choisissant n'importe quelle valeur pour $\widehat{\beta}_\lambda$, grâce au facteur multiplicatif \widehat{p}_λ .

Rappelons qu'on a alors calculé les termes principaux de la pénalité idéale

$$\text{pen}_{\text{id}}(m) = p_1(m) + p_2(m) - \delta(m) \quad \text{où} \quad \delta(m) = (P_n - P)\gamma(s_m^*), \quad (8)$$

le terme $\delta(m)$ étant d'espérance nulle.

Espérances conditionnelles. En vue d'obtenir une formule pour l'excès de risque et l'espérance de la pénalité idéale, commençons par calculer les espérances de p_1 et p_2 conditionnellement aux positions

$$\mathcal{P}_m = \mathcal{P}_m(D_n) := (\mathbb{1}_{X_i \in \lambda})_{1 \leq i \leq n, \lambda \in m}$$

des X_i dans les éléments de m . Pour tout λ , sachant \mathcal{P}_m , si $\widehat{p}_\lambda > 0$, $\widehat{\beta}_\lambda$ est la moyenne de $n\widehat{p}_\lambda$ variables aléatoires i.i.d. de moyenne β_λ et de variance σ_λ^2 , où

$$\begin{aligned} \sigma_\lambda^2 &:= \mathbb{E}_{(X,Y) \sim P} \left[(Y - \beta_\lambda)^2 \mid X \in \lambda \right] = \left(\sigma_\lambda^{(d)} \right)^2 + \left(\sigma_\lambda^{(a)} \right)^2 \\ \text{et} \quad \left(\sigma_\lambda^{(a)} \right)^2 &:= \mathbb{E}_{(X,Y) \sim P} \left[(Y - s^*(X))^2 \mid X \in \lambda \right] \\ &= \mathbb{E}_{(X,Y) \sim P} \left[(\sigma(X))^2 \mid X \in \lambda \right]. \end{aligned}$$

Par conséquent, pour tout $\lambda \in m$ tel que $\widehat{p}_\lambda > 0$,

$$\mathbb{E} \left[\left(\widehat{\beta}_\lambda - \beta_\lambda \right)^2 \mid \mathcal{P}_m \right] = \frac{\sigma_\lambda^2}{n\widehat{p}_\lambda},$$

et donc, en supposant $\min_{\lambda \in m} \widehat{p}_\lambda > 0$, (6) et (7) impliquent

$$\mathbb{E}[p_1(m) \mid \mathcal{P}_m] = \frac{1}{n} \sum_{\lambda \in m} \left(\frac{\sigma_\lambda^2 p_\lambda}{\widehat{p}_\lambda} \right) \quad (9)$$

$$\mathbb{E}[p_2(m) \mid \mathcal{P}_m] = \frac{1}{n} \sum_{\lambda \in m} \sigma_\lambda^2. \quad (10)$$

Espérances. Commençons par p_2 : d'après (10), on a

$$\mathbb{E}[p_2(m)] = \frac{1}{n} \sum_{\lambda \in m} \sigma_\lambda^2. \quad (11)$$

On voit avec (9) que le cas $\widehat{p}_\lambda = 0$ va poser problème pour parler d'espérance de $p_1(m)$. En effet, on ne peut jamais exclure cet événement (qui a une probabilité strictement positive sauf dans le cas trivial $p_\lambda = 1$), si bien qu'il va falloir choisir une convention pour $\widehat{\beta}_\lambda$ lorsque $\widehat{p}_\lambda = 0$. Cette convention est purement formelle (on ne l'utilisera jamais en pratique, considérant que $\widehat{\beta}_\lambda$ est indéfini lorsque $\widehat{p}_\lambda = 0$), mais indispensable pour une étude théorique rigoureuse. Notons enfin que ce problème ne se pose pas pour $p_2(m)$ dont l'espérance (11) est parfaitement définie en toute généralité. Nous ne parlerons plus dans ces notes de ce problème avant tout technique, renvoyant à [3, 4] pour donner un sens précis aux résultats énoncés dans la suite.

À l'aide de (9), on obtient l'expression suivante pour l'«espérance»¹ de $p_1(m)$:

$$\mathbb{E}[p_1(m)] = \frac{1}{n} \sum_{\lambda \in m} \left(\sigma_\lambda^2 \mathbb{E} \left[\frac{p_\lambda}{\widehat{p}_\lambda} \mid \widehat{p}_\lambda > 0 \right] \right) . \quad (12)$$

Pour rendre plus claire l'expression obtenue avec (12), on note

$$\delta_{n,p_\lambda} := \mathbb{E} \left[\frac{p_\lambda}{\widehat{p}_\lambda} \mid \widehat{p}_\lambda > 0 \right] - 1 ,$$

qui devrait intuitivement être petite dès que np_λ est assez grand, car alors $n\widehat{p}_\lambda$ se concentre autour de son espérance np_λ . Plus précisément, comme $n\widehat{p}_\lambda$ suit une loi binômiale de paramètres (n, p_λ) , le résultat suivant — extrait de [3, Lemme 3] — nous donne un contrôle non-asymptotique de δ_{n,p_λ} .

Lemme 1. *Soit $n \in \mathbb{N} \setminus \{0\}$, $p \in]0, 1]$, Z une variable aléatoire binômiale de paramètres (n, p) . On pose $\kappa_1 = 5.1$ et $\kappa_2 = 3.2$. Alors, si $np \geq 1$,*

$$1 - e^{-np} \leq \mathbb{E}[Z] \mathbb{E}[Z^{-1} \mid Z > 0] \leq \min \left\{ 1 + \frac{\kappa_1}{(np)^{1/4}} , \kappa_2 \right\} . \quad (13)$$

La preuve du Lemme 1 repose principalement sur l'inégalité de Bernstein [39, Proposition 2.9] appliquée à Z , voir [3].

En combinant (11), (12) et le Lemme 1, on obtient que p_1 et p_2 sont proches en espérance (clé de l'heuristique de pente, voir le deuxième cours et [6]) et l'expression suivante pour l'espérance de la pénalité idéale :

$$\begin{aligned} \mathbb{E}[\text{pen}_{\text{id}}(m)] &= \frac{1}{n} \sum_{\lambda \in m} \left[(2 + \delta_{n,p_\lambda}) (\sigma_\lambda)^2 \right] \\ &= \frac{1}{n} \sum_{\lambda \in m} \left[(2 + \delta_{n,p_\lambda}) \left((\sigma_\lambda^{(a)})^2 + (\sigma_\lambda^{(d)})^2 \right) \right] . \end{aligned} \quad (14)$$

De même, en combinant (5), (12) et le Lemme 1, on obtient une expression de l'excès de risque :

$$\mathbb{E}[\ell(s^*, \widehat{s}_m)] = \sum_{\lambda \in m} p_\lambda \left(\sigma_\lambda^{(d)} \right)^2 + \frac{1}{n} \sum_{\lambda \in m} \left[(1 + \delta_{n,p_\lambda}) \sigma_\lambda^2 \right] . \quad (15)$$

¹ou du moins une quantité déterministe autour de laquelle p_1 se concentre effectivement

Comparaison avec le cadre du deuxième cours. Il est intéressant de comparer l'expression (14) de l'espérance de la pénalité idéale avec l'expression $2\sigma^2 D_m/n$ obtenue au premier cours dans le cas homoscédastique avec un plan d'expérience déterministe. On peut relever deux différences principales :

- l'hétéroscédasticité (via les variations de $\sigma(\cdot)$) transforme $\sigma^2 D_m$ en $\sum_{\lambda \in m} \sigma_\lambda^{(a)2}$
- le caractère aléatoire du plan d'expérience ajoute les termes $\sigma_\lambda^{(d)2}$, qui sont le plus souvent négligeables car $\sigma_\lambda^{(d)2} \leq \|s^* - s_m^*\|_\infty^2 \rightarrow 0$ dès que $\sup_{\lambda \in m} \text{Leb}(\lambda) \rightarrow 0$ si s^* est un peu régulière, tandis que $\sigma_\lambda^{(a)2} \geq \inf_{x \in \lambda} (\sigma(x))^2$.

2. NÉCESSITÉ D'ESTIMER LA FORME DE LA PÉNALITÉ

L'objectif de cette section est de montrer que l'on ne peut se contenter d'une pénalité fonction de la dimension des modèles lorsque les données sont hétéroscédastiques.

2.1. Illustration sur un exemple. À titre d'illustration, considérons le cas $\mathcal{M}_n = \mathcal{M}_n^{(\text{reg}, 1/2)}$ avec $\mathcal{X} = [0, 1]$. Intuitivement, une telle collection de modèles laisse la possibilité de s'adapter à l'hétéroscédasticité, pour peu que l'on en tienne compte dans la pénalité. On peut commencer à formaliser cette intuition en calculant plus précisément l'espérance de la pénalité idéale. La quantité-clé apparaissant dans (14) est, pour $\lambda \subset [0, 1]$,

$$\begin{aligned} \sigma_\lambda^2 &= \mathbb{E} \left[(Y - \beta_\lambda)^2 \mid X \in \lambda \right] \\ &= \mathbb{E} \left[(s^*(X) - s_m^*(X))^2 \mid X \in \lambda \right] + \mathbb{E} \left[(\sigma(X))^2 \mid X \in \lambda \right] \\ &= \frac{1}{\text{Leb}(\lambda)} \int_\lambda (s^*(x) - s_m^*(x))^2 dx + \frac{1}{\text{Leb}(\lambda)} \int_\lambda (\sigma(x))^2 dx . \end{aligned}$$

Ainsi,

$$\begin{aligned} \mathbb{E}[p_2(m)] &= \frac{1}{n} \sum_{\lambda \in m} \sigma_\lambda^2 \\ &= \frac{2D_{m,1}}{n} \int_0^{1/2} \left[(s^*(X) - s_m^*(x))^2 + (\sigma(x))^2 \right] dx \\ &\quad + \frac{2D_{m,2}}{n} \int_{1/2}^1 \left[(s^*(X) - s_m^*(x))^2 + (\sigma(x))^2 \right] dx \\ &= \frac{2\sigma_a^2 D_{m,1}}{n} + \frac{2\sigma_b^2 D_{m,2}}{n} + R(m, n) \end{aligned}$$

$$\text{où } \sigma_a^2 := \int_0^{1/2} (\sigma(x))^2 dx \quad \sigma_b^2 := \int_{1/2}^1 (\sigma(x))^2 dx$$

$$\text{et } 0 \leq R(m, n) = \frac{2}{n} \left(D_{m,1} \int_0^{1/2} (s^*(X) - s_m^*(x))^2 dx \right.$$

$$\begin{aligned}
& + D_{m,2} \int_{1/2}^1 (s^*(X) - s_m^*(x))^2 dx \Big) \\
& \leq \frac{\ell(s^*, s_m^*)}{(\ln(n))^2},
\end{aligned}$$

car $D_{m,i} \leq n/(2(\ln(n))^2)$. On en déduit donc que

$$\mathbb{E}[\text{pen}_{\text{id}}(m)] \approx \frac{4}{n} [\sigma_a^2 D_{m,1} + \sigma_b^2 D_{m,2}],$$

qui est proportionnelle à $D_m = D_{m,1} + D_{m,2}$ si et seulement si

$$\sigma_a^2 = \int_0^{1/2} (\sigma(x))^2 dx = \sigma_b^2 = \int_{1/2}^1 (\sigma(x))^2 dx.$$

On va voir que lorsque cette condition est vérifiée, avec $\mathcal{M}_n = \mathcal{M}_n^{(\text{reg}, 1/2)}$, utiliser une pénalité fonction de la dimension est sous-optimal.

2.2. Caractérisation des pénalités fonction de la dimension. Il est très simple de caractériser l'ensemble des modèles pouvant être sélectionnés si l'on utilise une fonction de la dimension comme pénalité. En effet, pour tout $D \in \mathcal{D}_n = \{D_m \text{ t.q. } m \in \mathcal{M}_n\}$, on pose

$$\begin{aligned}
\mathcal{M}_{\text{dim}}(D) & := \operatorname{argmin}_{m \in \mathcal{M}_n \text{ t.q. } D_m = D} \{P_n \gamma(\hat{s}_m)\} \\
\text{et } \mathcal{M}_{\text{dim}} & := \bigcup_{D \in \mathcal{D}_n} \mathcal{M}_{\text{dim}}(D).
\end{aligned}$$

Alors, le lemme ci-dessous montre que toute procédure de choix de modèles par pénalisation avec une pénalité ne dépendant que de D_m sélectionne un modèle $\hat{m} \in \mathcal{M}_{\text{dim}}$.

Lemme 2. *Pour toute fonction $F : \mathcal{M}_n \mapsto \mathbb{R}$ et tout échantillon $(X_i, Y_i)_{1 \leq i \leq n}$,*

$$\operatorname{argmin}_{m \in \mathcal{M}_n} \{P_n \gamma(\hat{s}_m) + F(D_m)\} \subset \mathcal{M}_{\text{dim}}.$$

Preuve du Lemme 2. Soit $\hat{m}_F \in \operatorname{argmin}_{m \in \mathcal{M}_n} \{P_n \gamma(\hat{s}_m) + F(D_m)\}$. Alors, pour tout $m \in \mathcal{M}_n$,

$$P_n \gamma(\hat{s}_{\hat{m}_F}) + F(D_{\hat{m}_F}) \leq P_n \gamma(\hat{s}_m) + F(D_m). \quad (16)$$

En particulier, (16) est vraie pour tout $m \in \mathcal{M}_n$ tel que $D_m = D_{\hat{m}_F}$, pour lesquels $F(D_{\hat{m}_F}) = F(D_m)$. Par conséquent, (16) entraîne que $\hat{m}_F \in \mathcal{M}_{\text{dim}}(D_{\hat{m}_F})$, soit $\hat{m}_F \in \mathcal{M}_{\text{dim}}$. \square

Mentionnons ici que Breiman [13] avait remarqué que seul un petit nombre de modèles — appelés «RSS-extreme submodels» — peuvent être sélectionnés par des pénalités de la forme $F(D_m) = K D_m$ avec $K \geq 0$. Là où Breiman insistait sur les avantages de cette limitation du point de vue de la complexité algorithmique, nous allons voir que c'est précisément cette limitation qui empêche d'avoir un excès de risque équivalent à celui de l'oracle lorsque le niveau de bruit varie.

2.3. Sous-optimalité des pénalités fonction de la dimension. D'un point de vue théorique, en supposant pour simplifier que $\mathcal{M}_n = \mathcal{M}_n^{(\text{reg}, 1/2)}$, on peut prouver le résultat négatif suivant (extrait de [5]) pour les pénalités fonction de la dimension en présence d'hétéroscédasticité.

Théorème 1. *On suppose que $(X_1, Y_1), \dots, (X_n, Y_n) \in [0, 1] \times \mathbb{R}$ sont i.i.d. avec X_i uniforme sur $[0, 1]$ et $\forall i = 1, \dots, n, Y_i = X_i + \varepsilon_i$ avec $(\varepsilon_i)_{1 \leq i \leq n}$ indépendantes telles que $\mathbb{E}[\varepsilon_i | X_i] = 0$ et $\mathbb{E}[\varepsilon_i^2 | X_i] = (\sigma(X_i))^2$. On suppose de plus que $s^* \in \mathcal{C}^2$,*

$$\|\varepsilon_i\|_\infty \leq E < \infty, \quad \min \left\{ (\sigma_a)^2, (\sigma_b)^2 \right\} > 0 \quad \text{et} \quad (\sigma_a)^2 \neq (\sigma_b)^2$$

$$\text{où} \quad (\sigma_a)^2 := \int_0^{1/2} (\sigma(x))^2 dx \quad \text{et} \quad (\sigma_b)^2 := \int_{1/2}^1 (\sigma(x))^2 dx .$$

Soit $\mathcal{M}_n = \mathcal{M}_n^{(\text{reg}, 1/2)}$ la collection définie en Section 1.3, avec une dimension maximale $M_n = \lfloor n/(\ln(n))^2 \rfloor$. Alors, il existe des constantes $K_1, \mathbb{C}_1 > 0$ et un événement de probabilité au moins $1 - K_1 n^{-2}$ sur lequel, pour toute fonction $F : \mathcal{M}_n \mapsto \mathbb{R}$ et tout $\hat{m}_F \in \text{argmin}_{m \in \mathcal{M}_n} \{P_n \gamma(\hat{s}_m) + F(D_m)\}$,

$$\ell(s^*, \hat{s}_{\hat{m}_F}) \geq \mathbb{C}_1 \inf_{m \in \mathcal{M}_n} \{\ell(s^*, \hat{s}_m)\} \quad \text{avec} \quad \mathbb{C}_1 > 1 . \quad (17)$$

La constante \mathbb{C}_1 dépend uniquement de $(\sigma_a)^2 / (\sigma_b)^2$; la constante K_1 peut uniquement dépendre de $E, (\sigma_a)^2, (\sigma_b)^2, \|(s^*)'\|_\infty$ et $\|(s^*)''\|_\infty$.

Le Théorème 1 est prouvé dans [5]. Outre le calcul de l'espérance de $\text{pen}_{\text{id}}(m)$, il repose essentiellement sur un contrôle de l'erreur d'approximation $\ell(s^*, \hat{s}_m)$ en fonction de $(D_{m,1}, D_{m,2})$ et sur les inégalité de concentration suivantes pour les composantes de $\text{pen}_{\text{id}}(m)$:

Proposition 3. *Soit $\gamma > 0, B_n \geq 2$. Soit m une partition de \mathcal{X} telle que $\min_{\lambda \in m} \{n p_\lambda\} \geq B_n$ et $p_1(m), p_2(m)$ définies par (6) et (7). On suppose que*

(Ab) *Les données sont bornées : $\|Y_i\|_\infty \leq A < \infty$.*

(An) *Le niveau de bruit est uniformément minoré : $\sigma(X_i) \geq \sigma_{\min} > 0$ p.s.*

Alors, une constante absolue $L_1 > 0$ et une constante $L_2(A/\sigma_{\min}, \gamma)$ existent telles qu'avec probabilité au moins $1 - L_1 n^{-\gamma}$,

$$|p_1(m) - \mathbb{E}[p_1(m) | \mathcal{P}_m]| \leq L_2 \left[\frac{(\ln(n))^2}{\sqrt{D_m}} + D_m e^{-LB_n} \right] \mathbb{E}[p_2(m)] \quad (18)$$

$$|p_2(m) - \mathbb{E}[p_2(m) | \mathcal{P}_m]| \leq L_2 \frac{\ln(n)}{\sqrt{D_m}} \mathbb{E}[p_2(m)] . \quad (19)$$

De plus, sous la seule hypothèse **(Ab)**, l'inégalité de Bernstein [39, Proposition 2.9] implique : pour tout $x \geq 0$, avec probabilité au moins $1 - 2e^{-x}$,

$$\forall \theta \in]0, 1] , \quad |(P_n - P)(\gamma(s_m^*) - \gamma(s^*))| \leq \theta \ell(s^*, s_m^*) + \frac{6A^2 x}{\theta n} . \quad (20)$$

2.4. Pourquoi estimer la forme de la pénalité ? Le résultat montré par le Théorème 1 est très fort : toute pénalité de la forme $F(D_m)$, même avec

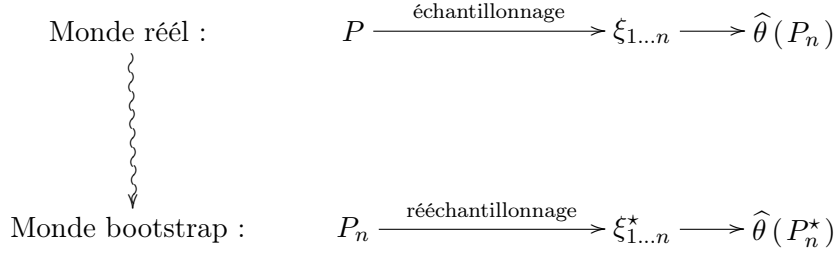


FIG. 1. L’heuristique de rééchantillonnage, selon Efron [19].
Schéma inspiré de [22, Figure 1].

une fonction F utilisant la connaissance complète de P et P_n , conduit à la perte d’un facteur multiplicatif $\mathbb{C}_1 > 1$ par rapport à l’oracle, en termes de perte relative. Si l’on veut une procédure satisfaisant une inégalité-oracle avec constante $(1 + o(1))$ dans un cadre hétéroscédastique, il est donc nécessaire d’estimer la forme de la pénalité.

Notons en revanche que si l’on se contente d’une inégalité-oracle avec constante multiplicative $\mathcal{O}(1)$, on peut utiliser une pénalité linéaire en D_m avec une constante suffisamment grande, par exemple $\text{pen}(m) = 2 \|\sigma\|_\infty^2 D_m n^{-1}$ [5, Proposition 2]. On risque toutefois le sur-apprentissage (et une explosion de la perte relative) si l’on sous-estime $\|\sigma\|_\infty^2$ [5, Proposition 3].

3. RÉÉCHANTILLONNAGE

Rééchantillonner, c’est utiliser un jeu de données (l’*échantillon*) pour construire un ou plusieurs nouveaux échantillons (les *rééchantillons*). Derrière une description aussi simple se cachent de nombreuses méthodes statistiques :

- le *sous-échantillonnage* consiste dans la sélection (aléatoire ou non) d’une partie des données initiales.
- le *bootstrap* (non-paramétrique) consiste dans la génération aléatoire d’un n -échantillon i.i.d., suivant la loi uniforme sur l’échantillon initialement observé. Autrement dit, on effectue n tirages avec remise dans les observations.

Dans ces deux cas, on peut — entre autres — estimer le biais ou la variance d’un estimateur en comparant le(s) rééchantillon(s) à l’échantillon initial, ou les rééchantillons entre eux. Pour une introduction plus complète au rééchantillonnage, nous renvoyons par exemple à [2, Section 1.1].

3.1. Heuristique d’Efron. L’heuristique qui sous-tend la plupart des utilisations du rééchantillonnage est le principe de «plug-in», décrit par Efron au sujet du bootstrap [19, 22]. Une autre description intuitive de cette heuristique se trouve dans l’introduction de [31]. Considérons la situation suivante : n observations $\xi_1, \dots, \xi_n \in \Xi$ ont été générées indépendamment avec une même loi P inconnue. On souhaite estimer un paramètre d’intérêt $\theta = t(P)$ (par exemple la moyenne, la variance, un quantile ; mais aussi des paramètres

plus complexes telles qu'une fonction de régression ou un prédicteur). Pour cela, on dispose d'un estimateur $\hat{\theta}(P_n) = \hat{\theta}(\xi_{1\dots n})$. La question qui se pose alors est : comment évaluer la qualité de cet estimateur ? C'est-à-dire, évaluer son biais, sa variance, construire un intervalle de confiance, ou — s'il s'agit d'un prédicteur — déterminer son erreur de prédiction.

L'heuristique proposée par Efron est de construire un «monde bootstrap», miroir du «monde réel» mais dans lequel aucune quantité n'est inconnue (voir Figure 1). La loi inconnue P est remplacée par une estimation \hat{P} . Nous nous limiterons ici au cadre «non-paramétrique» pour lequel \hat{P} est la mesure empirique $P_n = n^{-1} \sum_{i=1}^n \delta_{\xi_i}$. Le processus d'échantillonnage est remplacé par celui de rééchantillonnage. Dans le cas du bootstrap, les processus d'échantillonnage et de rééchantillonnage sont identiques : ξ_1^*, \dots, ξ_n^* sont i.i.d. de loi P_n . On note $P_n^* = n^{-1} \sum_{i=1}^n \delta_{\xi_i^*}$ leur mesure empirique. L'heuristique d'Efron se ramène donc à un principe de «plug-in» : on remplace P par son estimateur P_n , les autres quantités en découlant selon des processus inchangés.

La qualité d'un estimateur $\hat{\theta}(P_n)$ — de même que de nombreuses autres quantités — s'écrit sous la forme $F(P, P_n)$. Cette quantité est inaccessible car P est inconnue, mais son équivalent dans le monde bootstrap $F(P_n, P_n^*)$ ne dépend que des observations. On peut donc le calculer, au besoin en faisant une approximation de type Monte-Carlo (voir [31, Appendice II] et [23, Chapitre 23]). L'heuristique de rééchantillonnage se formalise donc ainsi, $\mathcal{L}(X)$ désignant la loi de la variable aléatoire X :

$$\mathcal{L}(F(P, P_n)) \approx \mathcal{L}(F(P_n, P_n^*) | P_n) .$$

3.2. Rééchantillonnage à poids échangeables. Dans l'heuristique ci-dessus, nous n'avons fait apparaître le rééchantillon que via sa mesure empirique

$$P_n^* = \frac{1}{n} \sum_{i=1}^n \delta_{\xi_i^*} = \frac{1}{n} \sum_{i=1}^n W_i \delta_{\xi_i}$$

où, pour tout i , on a noté W_i le nombre d'occurrences de ξ_i dans le rééchantillon $\xi_{1\dots n}^*$.

On peut ainsi reformuler et généraliser l'heuristique ci-dessus au «bootstrap à poids échangeables» (ou plus exactement rééchantillonnage à poids échangeables) [38, 44], en remplaçant P_n^* par

$$P_n^W := \frac{1}{n} \sum_{i=1}^n W_i \delta_{\xi_i}$$

où $W = (W_1, \dots, W_n) \in \mathbb{R}^n$ est un vecteur de poids aléatoire, échangeable² et indépendant de $\xi_{1\dots n}$. Notons que l'on ne suppose pas nécessairement que les W_i sont entiers, ni même positifs. Par contre, il est en général préférable d'avoir $\mathbb{E}[W_i] = 1$. Parmi les poids classiques, on trouve :

² $W \in \mathbb{R}^n$ est échangeable si et seulement si pour toute permutation σ de $\{1, \dots, n\}$, W a la même distribution que $(W_{\sigma(1)}, \dots, W_{\sigma(n)})$. Voir par exemple [1].

- *Efron* (m), $m \in \mathbb{N} \setminus \{0\}$: $mn^{-1}W$ suit une loi multinomiale de paramètres $(m; n^{-1}, \dots, n^{-1})$. Lorsque $m = n$, c'est le *bootstrap*. Lorsque $m < n$, c'est le «*m out of n*» *bootstrap*.
- *Random hold-out*³ (q), $q \in \{1, \dots, n-1\}$: $W_i = nq^{-1}\mathbb{1}_{i \in I}$ avec I un sous-ensemble aléatoire de $\{1, \dots, n\}$, choisi uniformément parmi les parties de cardinal q . Mise à part la contrainte sur la loi de I , il s'agit de poids de sous-échantillonnage (subsampling), également appelé «*bootstrap without replacement*». Lorsque $q = n-1$, on retrouve le *leave-one-out*.
- *Rademacher* (p), $p \in (0, 1)$: pW_1, \dots, pW_n sont i.i.d., de loi Bernoulli de paramètre p . Le nom Rademacher provient du cas $p = 1/2$ où les poids sont (à translation près) des variables Rademacher i.i.d. À la normalisation de P_n^W près, il s'agit de poids de sous-échantillonnage.

Mentionnons pour finir les poids *Poisson* (μ), $\mu > 0$: $\mu W_1, \dots, \mu W_n$ sont i.i.d. de loi Poisson de paramètre μ . Ces poids correspondent à un argument classique, dit de «*Poissonisation*», car on peut également les définir comme les poids Efron (M) avec $M \sim \mathcal{P}(\mu n)$. En prenant une taille de rééchantillon aléatoire, on peut ainsi rendre les poids W_i indépendants.

3.3. Arguments théoriques asymptotiques. Les études théoriques sur le rééchantillonnage se sont multipliées depuis les travaux d'Efron [19]. On peut citer par exemple [20, 23, 31, 46] sur le *bootstrap*, et [7, 25] sur le *bootstrap* à poids. Nous ne détaillerons ici qu'un seul résultat [50, Théorème 3.6.13] montrant que la version rééchantillonnée $\widehat{\mathbb{G}}_n^W$ du processus empirique recentré $\sqrt{n}(P_n - P)$ converge (à constante multiplicative près) vers le même processus Gaussien \mathbb{G} . Dans l'énoncé qui suit, on a omis des conditions de mesurabilité sur \mathcal{F} pour simplifier.

Théorème 2. *Soit \mathcal{F} une classe de Donsker de fonctions mesurables. Pour tout $n \in \mathbb{N}$, soit $(W_{n,1}, \dots, W_{n,n}) \in \mathbb{R}^n$ un vecteur aléatoire positif, échangeable, indépendant de $\xi_{1\dots n}$ tel que*

$$\sup_{n \in \mathbb{N}} \|W_{n,1} - \bar{W}_n\|_{2,1} < \infty \quad \text{avec} \quad \bar{W}_n = n^{-1} \sum_{i=1}^n W_{n,i}$$

$$\text{et} \quad \|Z\|_{2,1} := \int_0^\infty \sqrt{\mathbb{P}(|Z| > t)} dt ,$$

$$n^{-1/2} \mathbb{E} \left[\max_{1 \leq i \leq n} |W_{n,i} - \bar{W}_n| \right] \xrightarrow{(p)} 0 \quad \text{et} \quad n^{-1} \sum_{i=1}^n (W_{n,i} - \bar{W}_n)^2 \xrightarrow{(p)} c^2 > 0 .$$

Alors, lorsque n tend vers l'infini,

$$\sup_{h \in BL_1} \left| \mathbb{E}_W \left[h \left(\widehat{\mathbb{G}}_n^W \right) \right] - \mathbb{E} [h(c\mathbb{G})] \right| \xrightarrow{(p)} 0$$

$$\text{avec} \quad \widehat{\mathbb{G}}_n^W := \frac{1}{\sqrt{n}} \sum_{i=1}^n W_{n,i} (\delta_{\xi_i} - P_n) = \sqrt{n} (P_n^W - \bar{W}_n P_n) ,$$

³que l'on peut traduire par «*validation aléatoire*», à rapprocher de la validation croisée.

\mathbb{G} étant un processus gaussien de moyenne nulle et de fonction de covariance $\text{cov}(f, g) = P(fg) - P(f)P(g)$ et BL_1 l'ensemble des fonctions 1-lipschitziennes et bornées par 1 (pour la norme $\|\cdot\|_\infty$).

3.4. Usages du rééchantillonnage. De par sa formulation simple et très générale, le rééchantillonnage est désormais un outil statistique utilisé dans un grand nombre de domaines, voir notamment [53, 16]. À l'origine, l'objectif d'Efron se limitait à utiliser le bootstrap pour l'estimation du biais et de la variance d'un estimateur. Le bootstrap a ensuite été utilisé fructueusement pour :

- construire des intervalles de confiance [17],
- calculer des p -valeurs pour des statistiques de test [11, 10, 27],
- estimer une erreur de prédiction [51, 24, 41],
- faire de la sélection de modèles [21, 45],
- *etc.*

Notons que l'on peut utiliser d'autres types de rééchantillonnage, en particulier le sous-échantillonnage, pour construire des intervalles de confiance ou des tests [43], et pour bien d'autres applications encore.

Outre sa simplicité et sa généralité, le rééchantillonnage est souvent utilisé pour ses propriétés *stabilisatrices*. Ceci est particulièrement utile en classification, où nombre d'algorithmes sont très sensibles à une perturbation par un petit nombre de données. Une partie des observations étant absente dans chaque rééchantillon, on peut ainsi obtenir un algorithme dont la sortie varie très peu si l'on supprime un petit nombre de données. Citons ici le «bagging» (contraction de «bootstrap aggregating») [15] et les forêts aléatoires [14] entre autres algorithmes utilisant cette propriété du rééchantillonnage.

3.5. Limites du rééchantillonnage. La simplicité du rééchantillonnage le rend cependant facile à utiliser abusivement. Mentionnons ici quelques-unes de ses limites [37, 16] :

- l'estimation de queues de distributions plus lourdes que Gaussiennes [30].
- l'utilisation de poids échangeables correspond à l'hypothèse implicite que les données sont *échangeables*. Si ce n'est pas le cas, il faut alors faire très attention à la manière d'appliquer l'heuristique de rééchantillonnage. Par exemple, dans le cadre de la régression sur un plan d'expérience déterministe, les (X_i, Y_i) ne sont pas échangeables et l'on est amené à *rééchantillonner les résidus* lorsque les erreurs ε_i sont i.i.d. [23]. Le problème se corse en présence d'hétéroscédasticité, même s'il peut être contourné en choisissant bien les poids W [51].
- lorsque les données sont dépendantes, le bootstrap ne fonctionne pas en général [48], la manière la plus classique de résoudre le problème étant de rééchantillonner par blocs [42].

Avant d'utiliser le rééchantillonnage pour construire une pénalité, étudions de manière détaillée un estimateur par rééchantillonnage dans un cadre assez simple.

4. UN ESTIMATEUR DE LA VARIANCE PAR RÉÉCHANTILLONNAGE

Soient ξ_1, \dots, ξ_n des variables aléatoires i.i.d. de moyenne μ et de variance σ^2 . Leur variance s'écrit alors

$$\sigma^2 = \mathbb{E} \left[n \left(\frac{1}{n} \sum_{i=1}^n \xi_i - \mu \right)^2 \right] = n \mathbb{E} \left[(\mathbb{E}_{\xi \sim P_n} \xi - \mathbb{E}_{\xi \sim P} \xi)^2 \right] = n \mathbb{E} [F(P, P_n)] \quad , \quad (21)$$

si bien que l'on dispose de l'estimateur par rééchantillonnage

$$\widehat{\sigma_W^2} = n \mathbb{E}_W [F(P_n, P_n^W)] \quad .$$

4.1. Formule close. Calculons-le pour des poids W échangeables généraux :

$$\begin{aligned} \widehat{\sigma_W^2} &= n \mathbb{E}_W \left[\left(\frac{1}{\sum_{k=1}^n W_k} \sum_{i=1}^n (W_i \xi_i) - \frac{1}{n} \sum_{i=1}^n \xi_i \right)^2 \right] \\ &= \frac{1}{n} \mathbb{E}_W \left[\left(\sum_{i=1}^n \left[\left(\frac{nW_i}{\sum_{k=1}^n W_k} - 1 \right) \xi_i \right] \right)^2 \right] \\ &= \frac{1}{n} \mathbb{E}_W \left[\left(\sum_{i=1}^n \left[\left(\frac{nW_i}{\sum_{k=1}^n W_k} - 1 \right) (\xi_i - \mu) \right] \right)^2 \right] \end{aligned}$$

car

$$\sum_{i=1}^n \left(\frac{W_i}{\sum_{j=1}^n W_j} - 1 \right) = 0 \quad .$$

Ainsi,

$$\begin{aligned} \widehat{\sigma_W^2} &= \frac{1}{n} \mathbb{E}_W \left[\sum_{i,j} \left[\left(\frac{nW_i}{\sum_{k=1}^n W_k} - 1 \right) \left(\frac{nW_j}{\sum_{k=1}^n W_k} - 1 \right) (\xi_i - \mu) (\xi_j - \mu) \right] \right] \\ &= \frac{1}{n} \sum_{i=1}^n \left[R_V^{(W)} (\xi_i - \mu)^2 \right] + \frac{1}{n} \sum_{i \neq j} \left[R_C^{(W)} (\xi_i - \mu) (\xi_j - \mu) \right] \end{aligned}$$

en posant

$$\begin{aligned} R_V^{(W)} &:= \mathbb{E}_W \left[\left(\frac{nW_i}{\sum_{k=1}^n W_k} - 1 \right)^2 \right] \\ \text{et } R_C^{(W)} &:= \mathbb{E}_W \left[\left(\frac{nW_i}{\sum_{k=1}^n W_k} - 1 \right) \left(\frac{nW_j}{\sum_{k=1}^n W_k} - 1 \right) \right] \end{aligned}$$

pour des $i \neq j$ quelconques (ces quantités ne dépendent pas de (i, j) si $i \neq j$ car W est échangeable). Remarquons maintenant que

$$0 = \mathbb{E}_W \left[\left(\sum_{i=1}^n \left(\frac{nW_i}{\sum_{k=1}^n W_k} - 1 \right) \right)^2 \right] = nR_V^{(W)} + n(n-1)R_C^{(W)}$$

et donc, en supposant $n \geq 2$,

$$R_C^{(W)} = \frac{-1}{n-1} R_V^{(W)} .$$

Comme $R_V^{(W)} = 0$ lorsque $n = 1$, on en déduit que

$$\begin{aligned} \widehat{\sigma}_W^2 &= \frac{R_V^{(W)}}{n} \mathbf{1}_{n \geq 2} \left[\sum_{i=1}^n (\xi_i - \mu)^2 - \frac{1}{n-1} \sum_{i \neq j} ((\xi_i - \mu)(\xi_j - \mu)) \right] \quad (22) \\ &= \frac{R_V^{(W)}}{n} \mathbf{1}_{n \geq 2} \left[\sum_{i=1}^n \xi_i^2 - \frac{1}{n-1} \sum_{i \neq j} (\xi_i \xi_j) \right] , \end{aligned}$$

la dernière expression provenant de l'invariance de $\widehat{\sigma}_W^2$ par translation des données.

4.2. Comparaison avec l'estimateur classique. On suppose désormais que $n \geq 2$. Remarquons que l'estimateur sans biais classique de la variance s'écrit

$$\begin{aligned} \widehat{\sigma}^2 &= \frac{1}{n-1} \sum_{i=1}^n \left(\xi_i - \frac{1}{n} \sum_{k=1}^n \xi_k \right)^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n \left((\xi_i - \mu) - \frac{1}{n} \sum_{k=1}^n (\xi_k - \mu) \right)^2 \\ &= \frac{n}{n-1} \left[\frac{1}{n} \sum_{i=1}^n (\xi_i - \mu)^2 - \left(\frac{1}{n} \sum_{i=1}^n (\xi_i - \mu) \right)^2 \right] \\ &= \frac{1}{n} \left[\sum_{i=1}^n (\xi_i - \mu)^2 - \frac{1}{n-1} \sum_{i \neq j} ((\xi_i - \mu)(\xi_j - \mu)) \right] . \quad (23) \end{aligned}$$

Ainsi,

$$\widehat{\sigma}_W^2 = R_V^{(W)} \widehat{\sigma}^2 ,$$

d'où l'on déduit notamment que

$$\mathbb{E} \left[\widehat{\sigma}_W^2 \right] = R_V^{(W)} \sigma^2 . \quad (24)$$

4.3. Rééchantillonnage et structure. Une explication de la relation (24) est que l'estimateur par rééchantillonnage de $\mathbb{E}[F(P, P_n)]$ respecte la *structure* de $F(P, P_n)$. En effet, lorsque $F(P, P_n)$ est définie par (21), elle possède

les trois propriétés suivantes :

échangeabilité : $F(P, P_n)$ ne dépend pas de l'ordre des données (25)

invariance par translation de la loi de ξ d'une constante $c \in \mathbb{R}$ (26)

homogénéité : si $\xi - \mu$ est multipliée par $\lambda > 0$, (27)

alors $F(P, P_n)$ est multipliée par λ^2 .

De plus, $F(P, P_n)$ ne dépend de P que via une quantité de la forme $\mathbb{E}_{\xi \sim P} Q(\xi)$ pour un polynôme Q , et $F(P, P_n)$ est un polynôme en $\mathbb{E}_{\xi \sim P} Q(\xi)$ et ξ_1, \dots, ξ_n . Par conséquent, lorsque W est échangeable,

$$\widehat{\sigma_W^2} = \mathbb{E}_W [F(P_n, P_n^W)]$$

possède automatiquement les mêmes propriétés structurelles.

Comme remarqué dans [2, Section 8.2], ces conditions sont suffisamment fortes pour contraindre $\widehat{\sigma_W^2}$ à satisfaire (24) pour une constante $R_V^{(W)}$ dès lors que W est échangeable. Plus précisément, on a le lemme suivant, prouvé initialement dans [2, Lemme 8.2].

Lemme 4. Soit $n \in \mathbb{N} \setminus \{0\}$ et $f : z = (z_i)_{1 \leq i \leq n} \in \mathbb{R}^n \mapsto \mathbb{R}$ vérifiant :

(25) f est échangeable, c'est-à-dire, pour toute permutation σ de $\{1, \dots, n\}$ et tout $z \in \mathbb{R}^n$, $f((z_{\sigma(i)})_{1 \leq i \leq n}) = f(z)$.

(26) pour tout $c \in \mathbb{R}$ et $z \in \mathbb{R}^n$, $f((z_i + c)_{1 \leq i \leq n}) = f(z)$.

(27) $f(z)$ est un polynôme en z_1, \dots, z_n et pour tout $\lambda \in \mathbb{R}$ et $z \in \mathbb{R}^n$, $f((\lambda z_i)_{1 \leq i \leq n}) = \lambda^2 f(z)$.

Alors, il existe $\alpha \in \mathbb{R}$ tel que pour tout $z \in \mathbb{R}^n$,

$$f(z) = \alpha \left(\sum_{i=1}^n z_i^2 - \frac{1}{n} \left(\sum_{1 \leq i \leq n} z_i \right)^2 \right) .$$

Réciproquement, une fonction de cette forme possède les propriétés (25)–(27).

Preuve du Lemme 4. D'après (27), il existe $(a_{i,j})_{1 \leq i,j \leq n} \in \mathbb{R}^{n^2}$ tel que pour tout $z \in \mathbb{R}^n$,

$$f(z) = \sum_{1 \leq i,j \leq n} a_{i,j} z_i z_j .$$

D'après (25), $a_{i,j}$ peut uniquement dépendre de $\mathbf{1}_{i=j}$, c'est-à-dire, pour tout $z \in \mathbb{R}^n$,

$$f(z) = \alpha \sum_{i=1}^n z_i^2 + 2\beta \sum_{1 \leq i < j \leq n} z_i z_j = (\alpha - \beta) \sum_{i=1}^n z_i^2 + \beta \left(\sum_{i=1}^n z_i \right)^2$$

avec $a_{1,1} = \alpha$ et $a_{1,2} = \beta$. En utilisant maintenant (26), on a pour tout $c \in \mathbb{R}$ et $z \in \mathbb{R}^n$,

$$f(z + c) = f(z) + (\alpha + (n-1)\beta) \left(2c \sum_{i=1}^n z_i + nc^2 \right) .$$

On a donc $\alpha + (n - 1)\beta = 0$, d'où le résultat. La réciproque se vérifie simplement. \square

4.4. Rééchantillonnage et concentration. Il est intéressant de comparer la variance de

$$Z_1 = F(P, P_n) = n \left(\frac{1}{n} \sum_{i=1}^n \xi_i - \mu \right)^2 = \frac{1}{n} \left(\sum_{i=1}^n (\xi_i - \mu) \right)^2$$

à celle de l'estimateur par rééchantillonnage de son espérance $\widehat{\sigma_W^2}$.

Le Lemme 5 ci-dessous (avec $a = 0$ et $b = 1/n$) montre que

$$\text{var}(Z_1) = 2\sigma^4 + \frac{\mathbb{E}[(\xi_1 - \mu)^4] - 3\sigma^4}{n} . \quad (28)$$

De même, le Lemme 5 (avec $a = 1/(n - 1)$ et $b = -1/(n(n - 1))$) montre que

$$\text{var} \left(\frac{1}{R_V^{(W)}} \widehat{\sigma_W^2} \right) = \frac{1}{n} \left(\mathbb{E}[(\xi_1 - \mu)^4] - \sigma^4 \right) + \frac{2}{n(n - 1)} \sigma^4 . \quad (29)$$

La comparaison de (28) et (29) met en lumière un phénomène de sur-concentration, classique pour les estimateurs d'espérances par rééchantillonnage : la variance de $\widehat{\sigma_W^2}$ est de l'ordre de n^{-1} alors que la quantité-cible $F(P, P_n)$ a une variance de l'ordre de 1.

Pour finir, énonçons et prouvons le Lemme 5.

Lemme 5. Soit ξ_1, \dots, ξ_n des v.a.i.i.d. centrées de variance σ^2 , et l'on pose $\kappa^4 = \mathbb{E}[\xi_1^4]$. Soient $a, b \in \mathbb{R}$ et

$$Z = a \sum_{i=1}^n \xi_i^2 + b \left(\sum_{i=1}^n \xi_i \right)^2 = (a + b) \sum_{i=1}^n \xi_i^2 + b \sum_{i \neq j} \xi_i \xi_j .$$

Alors,

$$\text{var}(Z) = n(a + b)^2 (\kappa^4 - \sigma^4) + 2n(n - 1)b^2 \sigma^4 .$$

Preuve du Lemme 5. D'une part,

$$\mathbb{E}[Z] = (a + b)n\sigma^2 .$$

D'autre part,

$$\begin{aligned} \mathbb{E}[Z^2] &= a^2 \mathbb{E} \left[\left(\sum_{i=1}^n \xi_i^2 \right)^2 \right] + b^2 \mathbb{E} \left[\left(\sum_{i=1}^n \xi_i \right)^4 \right] + 2ab \sum_{i,j,k} \mathbb{E}[\xi_i^2 \xi_j \xi_k] \\ &= (a + b)^2 \sum_{i,j} \mathbb{E}[\xi_i^2 \xi_j^2] + 2b^2 \sum_{i \neq j} \mathbb{E}[\xi_i^2 \xi_j^2] \\ &= (a + b)^2 (n\kappa^4 + n(n - 1)\sigma^4) + 2b^2 n(n - 1)\sigma^4 . \end{aligned}$$

On en déduit le résultat en retirant $(\mathbb{E}Z)^2$ à $\mathbb{E}[Z^2]$. \square

$\mathcal{L}(W)$	Efr(m)	Rad(p)	Rho(q)	Loo
$R_V^{(W)}$	$\frac{n-1}{m}$	$\frac{1+\delta_{n,p}}{p} - 1$	$\frac{n}{q} - 1$	$\frac{1}{n-1}$

TAB. 1. Valeur de $R_V^{(W)}$, défini par (30), pour différents ré-échantillonnages.

4.5. Calcul de la constante multiplicative. Pour finir cette section, calculons pour les exemples classiques de poids la valeur de la constante multiplicative

$$R_V^{(W)} := \mathbb{E}_W \left[\left(\frac{nW_i}{\sum_{k=1}^n W_k} - 1 \right)^2 \right] \quad (30)$$

par laquelle il faut normaliser $\widehat{\sigma_W^2}$ pour en faire un estimateur sans biais de σ^2 . La Table 1 synthétise les résultats obtenus.

Démonstration. Pour les poids Efron(m) et Random hold-out (q), on a

$$\sum_{i=1}^n W_i = n \quad \text{p.s.} \quad \text{et} \quad \mathbb{E}[W_i] = 1 \quad ,$$

si bien que $R_V^{(W)} = \text{var}(W_1)$. On en déduit :

– pour Efron(m) (abrégé Efr(m)) : $W_1 \sim (n/m) \times \mathcal{B}(m, 1/n)$ et donc

$$R_V^{(W)} = \text{var}(W_1) = \frac{n^2}{m^2} \times m \times \frac{1}{n} \left(1 - \frac{1}{n} \right) = \frac{n-1}{m} \quad .$$

– pour Random hold-out (q) (abrégé Rho(q)) : $W_1 \sim (n/q) \times \mathcal{B}(q/n)$ et donc

$$R_V^{(W)} = \text{var}(W_1) = \frac{n^2}{q^2} \times \frac{q}{n} \left(1 - \frac{q}{n} \right) = \frac{n}{q} - 1 \quad .$$

En prenant $q = n - 1$, on obtient la valeur de $R_V^{(W)}$ pour les poids Leave-one-out (abrégé Loo).

Pour les poids Rademacher (p) (abrégé Rad(p)), rappelons que les W_i sont i.i.d. de loi $(1/p) \times \mathcal{B}(p)$. La somme des W_i n'étant pas déterministe, nous allons commencer par raisonner conditionnellement à $\overline{W} = n^{-1} \sum_{i=1}^n W_i$. Comme W_i ne prend que deux valeurs 0 ou $1/p$, sachant \overline{W} , $W_i \sim (1/p) \times \mathcal{B}(p\mathbb{E}[W_i | \overline{W}])$. Or, les W_i étant échangeables, $\mathbb{E}[W_i | \overline{W}]$ ne dépend pas de i , si bien que

$$\mathbb{E}[W_i | \overline{W}] = \frac{1}{n} \sum_{k=1}^n \mathbb{E}[W_k | \overline{W}] = \mathbb{E}[\overline{W} | \overline{W}] = \overline{W} \quad .$$

On en déduit que

$$\mathbb{E} \left[(W_i - \overline{W})^2 \mid \overline{W} \right] = \text{var}(W_i | \overline{W}) = \frac{1}{p^2} p \overline{W} (1 - p \overline{W}) = \frac{\overline{W}}{p} (1 - p \overline{W}) \quad ,$$

et donc

$$R_V^{(W)} = \mathbb{E} \left[\frac{(W_i - \bar{W})^2}{\bar{W}^2} \right] = \mathbb{E} \left[\frac{1}{p\bar{W}} (1 - p\bar{W}) \right] = \mathbb{E} \left[\frac{1}{p\bar{W}} \right] - 1 .$$

Or, $np\bar{W} \sim \mathcal{B}(n, p)$, si bien que l'on peut utiliser le Lemme 1 :

$$\mathbb{E} \left[\frac{1}{\bar{W}} \right] = 1 + \delta_{n,p}$$

avec $\delta_{n,p} \rightarrow 0$ lorsque $np \rightarrow \infty$, si bien que

$$R_V^{(W)} = \frac{1 + \delta_{n,p}}{p} - 1 \approx \frac{1}{p} - 1 .$$

□

5. PÉNALITÉS PAR RÉÉCHANTILLONNAGE

Nous avons vu (notamment avec le Lemme 2 du premier cours) qu'estimer sans biais l'espérance de la pénalité idéale

$$\text{pen}_{\text{id}}(m) = (P - P_n)\gamma(\hat{s}_m(P_n))$$

est une bonne stratégie pour obtenir une procédure de choix de modèles par pénalisation optimale. Comme la pénalité idéale s'écrit comme une fonction de (P, P_n) (et ceci en supposant seulement que les estimateurs \hat{s}_m ne sont pas fonctions de l'ordre des données, mais seulement des fonctions de P_n), l'heuristique de rééchantillonnage (Section 3.1) nous conduit directement à l'estimateur suivant de l'espérance de $\text{pen}_{\text{id}}(m)$, que nous appellerons pénalité par rééchantillonnage :

$$\text{pen}_W(m) := C_W \mathbb{E}_W [(P_n - P_n^W)\gamma(\hat{s}_m(P_n^W))] \quad (31)$$

La pénalité (31) a d'abord été proposée par Efron dans le cas du bootstrap [21], avec un premier résultat théorique sur son optimalité asymptotique dans [47]. La même pénalité avec des poids Efron(m) a été proposée par Shao [45] dans un objectif d'identification (voir Section 3.2 du premier cours), en prenant $m \ll n$. La définition (31) avec des poids échangeables généraux provient de [4].

Notons la présence d'une constante C_W devant la pénalité (31). Elle joue le même rôle que le facteur $1/R_V^{(W)}$ nécessaire pour l'estimation de la variance en Section 4. Nous proposerons une valeur pour C_W pour les poids les plus classiques en Section 6, dans le cas des régressogrammes, mais sans garantie aucune sur leur validité dans un cadre général. Il nous semble raisonnable de proposer d'estimer la forme de la pénalité seulement à l'aide de (31), et de calibrer la constante multiplicative à l'aide de l'heuristique de pente lorsque c'est possible (voir le deuxième cours).

D'autres pénalités par rééchantillonnage ont été proposées, en particulier les complexités de Rademacher (globales [33, 8] ou locales [9, 32]) en classification, les premières ayant été généralisées avec les pénalités globales par

rééchantillonnage échangeables en classification [25]. Nous ne parlerons ici (brièvement) que des complexités globales ; voir aussi [12, 39] au sujet des complexités globales ou locales en classification.

Les complexités de Rademacher (globales) reposent sur la majoration suivante de la pénalité idéale lorsque $\widehat{s}_m \in S_m$ p.s. :

$$\text{pen}_{\text{id}}(m) = (P - P_n)\gamma(\widehat{s}_m(P_n)) \leq \sup_{t \in S_m} \{(P - P_n)\gamma(t)\} =: \text{pen}_{\text{id,g}}(m) .$$

Une telle majoration peut paraître large, mais elle reste raisonnable pour obtenir des vitesses standard en classification avec un contraste borné. Son intérêt est que pour obtenir une inégalité-oracle, nous avons vu avec le Lemme 2 du premier cours qu'il est plus important d'avoir la minoration $\text{pen}(m) \geq \text{pen}_{\text{id}}(m)$ que d'avoir une majoration de $\text{pen}(m)$: au pire, le membre de droite de l'inégalité-oracle est augmenté de $\text{pen}(m) - \text{pen}_{\text{id}}(m)$, mais surpénaliser ne risque pas de conduire au sur-apprentissage. Ainsi, utiliser une estimation de $\text{pen}_{\text{id,g}}(m)$ (ou de son espérance) comme pénalité peut sembler raisonnable.

À l'aide d'un argument de symétrisation, on a obtenu la complexité de Rademacher globale

$$\mathbb{E} \left[\sup_{t \in S_m} \left\{ \frac{1}{n} \sum_{i=1}^n (\varepsilon_i \gamma(t; \xi_i)) \right\} \middle| P_n \right] \quad \text{avec} \quad \varepsilon_1, \dots, \varepsilon_n \text{ i.i.d. } \sim \mathcal{U}(\{-1, +1\}) \quad (32)$$

indépendantes de P_n . À constante multiplicative près (comprise entre un et deux), on peut montrer que (32) est proche de $\text{pen}_{\text{id,g}}(m)$, d'où une inégalité-oracle (fine si $\text{pen}_{\text{id,g}}$ n'est pas trop grande par rapport à pen_{id}). Mais si la principale justification théorique de (32) est un argument de symétrisation, il est intéressant de noter que cette pénalité est l'estimation par rééchantillonnage (avec des poids Rademacher (1/2), d'où leur nom) de $\text{pen}_{\text{id,g}}$. En effet,

$$\begin{aligned} & \mathbb{E}_W \left[\sup_{t \in S_m} \{(P_n - P_n^W)\gamma(t)\} \middle| P_n \right] \\ &= \mathbb{E}_W \left[\sup_{t \in S_m} \left\{ \frac{1}{n} \sum_{i=1}^n ((1 - W_i)\gamma(t; \xi_i)) \right\} \middle| P_n \right] \end{aligned} \quad (33)$$

coïncide avec (32) car $1 - W_i \sim \varepsilon_i$ pour les poids Rademacher (1/2). En particulier, il est naturel de généraliser (32) à des poids généraux via (33), ce qui a été fait dans [25].

6. GARANTIES THÉORIQUES POUR LES RÉGRESSOGRAMMES

Cette section rassemble des résultats issus de [4] qui donnent des garanties théoriques pour les pénalités par rééchantillonnage échangeables (31) pour la sélection de régressogrammes en régression hétéroscédastique sur un plan d'expérience aléatoire (voir Section 1).

6.1. Calcul de la pénalité par rééchantillonnage. Commençons par calculer explicitement la pénalité (31) lorsque \hat{s}_m est un régressogramme associé à une partition finie m de \mathcal{X} . En appliquant directement la formule (31) au cadre des régressogrammes, on obtient la pénalité

$$C_W \mathbb{E}_W [P_n \gamma(\hat{s}_m^W) - P_n^W \gamma(\hat{s}_m^W)] \quad , \quad (34)$$

$$\text{où } \hat{s}_m^W := \operatorname{argmin}_{t \in S_m} P_n^W \gamma(t) = \sum_{\lambda \in m} \hat{\beta}_\lambda^W \mathbf{1}_\lambda, \quad \hat{\beta}_\lambda^W := \frac{1}{n \hat{p}_\lambda^W} \sum_{X_i \in \lambda} W_i Y_i \quad ,$$

$$\hat{p}_\lambda^W := P_n^W(X \in \lambda) = \hat{p}_\lambda \hat{W}_\lambda \quad \text{et} \quad \hat{W}_\lambda := \frac{1}{n \hat{p}_\lambda} \sum_{X_i \in \lambda} W_i \quad .$$

Convention. De la même manière que l'espérance de $p_1(m)$ n'est pas bien définie à cause de l'événement $\min_{\lambda \in m} \hat{p}_\lambda = 0$ (voir (12)), l'événement $\min_{\lambda \in m} \hat{p}_\lambda^W = 0$ pose problème dans la définition (34) et a en général une probabilité non-nulle. Commençons par séparer la pénalité par rééchantillonnage en trois termes, à la suite de la décomposition (8) de la pénalité idéale :

$$\mathbb{E}_W [P_n \gamma(\hat{s}_m^W) - P_n^W \gamma(\hat{s}_m^W)] = \hat{p}_0(m) + \hat{p}_1(m) + \hat{p}_2(m)$$

avec

$$\hat{p}_0(m) := \mathbb{E}_W [(P_n - P_n^W) \gamma(\hat{s}_m)] = \frac{1}{n} \sum_{i=1}^n (\mathbb{E}_W [1 - W_i] \gamma(\hat{s}_m; (X_i, Y_i))) = 0$$

car $\mathbb{E}_W[W_i] = 1$ pour tout i ,

$$\hat{p}_1(m) := \sum_{\lambda \in m} \mathbb{E}_W \left[\hat{p}_\lambda \left(\hat{\beta}_\lambda^W - \hat{\beta}_\lambda \right)^2 \right]$$

$$\text{et } \hat{p}_2(m) := \sum_{\lambda \in m} \mathbb{E}_W \left[\hat{p}_\lambda^W \left(\hat{\beta}_\lambda^W - \hat{\beta}_\lambda \right)^2 \right] \quad .$$

Avec la convention $\hat{p}_\lambda^W (\hat{\beta}_\lambda^W - \hat{\beta}_\lambda)^2 = 0$ lorsque $\hat{p}_\lambda^W = 0$, $\hat{p}_2(m)$ est bien définie car $\hat{\beta}_\lambda^W$ est bien défini lorsque $\hat{p}_\lambda^W > 0$. Il reste à définir $\hat{p}_1(m)$. Notre proposition est de remplacer l'espérance sur W par une espérance conditionnellement à l'événement $\hat{W}_\lambda > 0$, *séparément pour chaque* $m \in \mathcal{M}_n$ et $\lambda \in m$, ce qui garantit que la grande majorité des vecteurs de poids restent admissible; voir [2, Section 8.1] pour une justification de ce choix dans un cadre plus général. En résumé, (31) est remplacée par

$$C_W \sum_{\lambda \in m} \left(\mathbb{E}_W \left[\hat{p}_\lambda \left(\hat{\beta}_\lambda^W - \hat{\beta}_\lambda \right)^2 \mid \hat{W}_\lambda > 0 \right] + \mathbb{E}_W \left[\hat{p}_\lambda^W \left(\hat{\beta}_\lambda^W - \hat{\beta}_\lambda \right)^2 \right] \right) \quad . \quad (35)$$

Formule close. Nous pouvons désormais énoncer une formule close pour la pénalité par rééchantillonnage, correctement définie par (35).

Lemme 6. *Soit m une partition de \mathcal{X} et \hat{s}_m le régressogramme associé. Soit $W \in [0, \infty[^n$ un vecteur aléatoire échangeable indépendant de $(X_i, Y_i)_{1 \leq i \leq n}$.*

Soit $\text{pen}_W(m)$ définie par (35). On suppose que $\min_{\lambda \in m} \{n\hat{p}_\lambda\} \geq 1$. Alors,

$$\text{pen}_W(m) = \frac{C_W}{n} \sum_{\lambda \in m} (R_{1,W}(n, \hat{p}_\lambda) + R_{2,W}(n, \hat{p}_\lambda)) \frac{n\hat{p}_\lambda S_{\lambda,2} - S_{\lambda,1}^2}{n\hat{p}_\lambda(n\hat{p}_\lambda - 1)} \mathbb{1}_{n\hat{p}_\lambda \geq 2} \quad (36)$$

$$\text{avec } S_{\lambda,1} := \sum_{X_i \in \lambda} (Y_i - \beta_\lambda) \quad S_{\lambda,2} := \sum_{X_i \in \lambda} (Y_i - \beta_\lambda)^2 \quad (37)$$

$$R_{1,W}(n, \hat{p}_\lambda) := \mathbb{E} \left[\frac{(W_1 - \widehat{W}_\lambda)^2}{\widehat{W}_\lambda^2} \middle| X_1 \in \lambda, \widehat{W}_\lambda > 0 \right] \quad (38)$$

$$\text{et } R_{2,W}(n, \hat{p}_\lambda) := \mathbb{E} \left[\frac{(W_1 - \widehat{W}_\lambda)^2}{\widehat{W}_\lambda} \middle| X_1 \in \lambda \right]. \quad (39)$$

Preuve du Lemme 6. La pénalité idéale s'écrit sous la forme

$$\text{pen}_{\text{id}}(m) = \sum_{\lambda \in m} F_\lambda(\mathcal{L}((X, Y) | X \in \lambda), (X_i, Y_i)_{X_i \in \lambda}).$$

Le rééchantillonnage conservant cette structure, il n'est pas surprenant que l'on puisse décomposer de même la pénalité par rééchantillonnage :

$$\text{pen}_W(m) = C_W (\hat{p}_1(m) + \hat{p}_2(m))$$

$$\begin{aligned} \text{avec } \hat{p}_1(m) &:= \sum_{\lambda \in m} \mathbb{E}_W \left[\hat{p}_\lambda \left(\hat{\beta}_\lambda^W - \beta_\lambda \right)^2 \middle| \widehat{W}_\lambda > 0 \right] \\ &= \frac{1}{n} \sum_{\lambda \in m} \mathbb{E}_W \left[\hat{v}_\lambda^W \left(\widehat{W}_\lambda \right) \middle| \widehat{W}_\lambda > 0 \right] \end{aligned}$$

$$\hat{p}_2(m) := \sum_{\lambda \in m} \mathbb{E}_W \left[\hat{p}_\lambda^W \left(\hat{\beta}_\lambda^W - \beta_\lambda \right)^2 \right] = \frac{1}{n} \sum_{\lambda \in m} \mathbb{E}_W \left[\widehat{W}_\lambda \hat{v}_\lambda^W \left(\widehat{W}_\lambda \right) \right]$$

$$\hat{v}_\lambda^W \left(\widehat{W}_\lambda \right) = \hat{v}_\lambda^W \left(\widehat{W}_\lambda, (X_i, Y_i)_{X_i \in \lambda} \right) := n\hat{p}_\lambda \mathbb{E}_W \left[\left(\hat{\beta}_\lambda^W - \beta_\lambda \right)^2 \middle| \widehat{W}_\lambda \right].$$

Notons que le conditionnement suivant la valeur de \widehat{W}_λ permet de faire apparaître la quantité que nous avons calculée à la Section 4. En effet, sachant \widehat{W}_λ , les poids $(W_i)_{X_i \in \lambda}$ sont échangeables et indépendants de $(X_i, Y_i)_{X_i \in \lambda}$, si bien que $\hat{v}_\lambda^W \left(\widehat{W}_\lambda \right)$ est l'estimateur par rééchantillonnage échangeable de la variance σ_λ^2 de Y sachant $X \in \lambda$. Ceci permet non seulement de comprendre pourquoi la pénalité par rééchantillonnage estime bien l'espérance de la pénalité idéale (14) malgré l'hétéroscédasticité, mais également d'obtenir que

$$\hat{v}_\lambda^W \left(\widehat{W}_\lambda \right) = \frac{R_V^{(W)}(\widehat{W}_\lambda, n\hat{p}_\lambda)}{n\hat{p}_\lambda - 1} \left[S_{\lambda,2} - \frac{1}{n\hat{p}_\lambda} S_{\lambda,1}^2 \right] \mathbb{1}_{n\hat{p}_\lambda \geq 2}$$

$$\text{avec } R_V^{(W)}(\widehat{W}_\lambda, n\hat{p}_\lambda) := \mathbb{E}_W \left[\left(\frac{W_1}{\widehat{W}_\lambda} - 1 \right)^2 \middle| \widehat{W}_\lambda, X_1 \in \lambda \right].$$

On en déduit le résultat annoncé. \square

$\mathcal{L}(W)$	Efr(m)	Rad(p)	Poi(μ)	Rho(q)	Loo
$R_{2,W}(n, \widehat{p}_\lambda)$	$\frac{n}{m} \left(1 - \frac{1}{n\widehat{p}_\lambda}\right)$	$\frac{1}{p} - 1$	$\frac{1}{\mu} \left(1 - \frac{1}{n\widehat{p}_\lambda}\right)$	$\frac{n}{q} - 1$	$\frac{1}{n-1}$
$C_{W,\infty}$	m/n	$p/(1-p)$	μ	$q/(n-q)$	$n-1$

TAB. 2. $R_{2,W}(n, \widehat{p}_\lambda)$ et $C_{W,\infty}$ pour quelques poids de rééchantillonnage.

Constantes multiplicatives pour des poids classiques. Pour finir de calculer complètement la pénalité par rééchantillonnage, il faut préciser les valeurs de $R_{1,W}(n, \widehat{p}_\lambda)$ et $R_{2,W}(n, \widehat{p}_\lambda)$. Pour les poids considérés en Section 4.5, on peut montrer que $R_{1,W}(n, \widehat{p}_\lambda) \approx R_{2,W}(n, \widehat{p}_\lambda)$ lorsque $n\widehat{p}_\lambda$ est assez grand, et l'on obtient les formules de la Table 2 pour $R_{2,W}(n, \widehat{p}_\lambda)$ [4, Lemme 17].

Espérances. Comme

$$\mathbb{E}[Y_i - \beta_\lambda \mid X_i \in \lambda] = 0 \quad \text{et} \quad \mathbb{E}\left[(Y_i - \beta_\lambda)^2 \mid X_i \in \lambda\right] = \sigma_\lambda^2,$$

on déduit simplement de (36) que

$$\mathbb{E}[\text{pen}_W(m) \mid \mathcal{P}_m] = \frac{C_W}{n} \sum_{\lambda \in m} \left[(R_{1,W}(n, \widehat{p}_\lambda) + R_{2,W}(n, \widehat{p}_\lambda)) \sigma_\lambda^2 \mathbf{1}_{n\widehat{p}_\lambda \geq 2} \right], \quad (40)$$

à comparer avec (9) et (10). En particulier, la pénalité par rééchantillonnage s'adapte bien à l'hétéroscédasticité des données dans le cas des régressogrammes.

En supposant $\min_{\lambda \in m} \{n\widehat{p}_\lambda\}$ suffisamment grand, compte-tenu de l'heuristique $R_{1,W}(n, \widehat{p}_\lambda) \approx R_{2,W}(n, \widehat{p}_\lambda)$, on en déduit que le bon choix (asymptotique) pour la constante C_W , celui qui conduit à une estimation sans biais de $\mathbb{E}[\text{pen}_{\text{id}}(m)]$, est de prendre

$$C_W = C_{W,\infty} := \lim_{n\widehat{p}_\lambda \rightarrow \infty} \frac{1}{R_{2,W}(n, \widehat{p}_\lambda)},$$

dont les valeurs⁴ sont reportées en Table 2 pour les poids les plus classiques. Il est remarquable que les valeurs classiques des paramètres des poids (Efr(n), Rad($1/2$), Poi(1), Rho($n/2$)) conduisent tous à une constante $C_{W,\infty} = 1$, essentiellement car ce sont les valeurs des paramètres pour lesquelles $\text{var}(W_1) \approx 1$.

Plus précisément, on peut montrer pour ces poids [4, Section 4] que

$$\mathbb{E}[\text{pen}_W(m)] = \frac{C_W}{C_{W,\infty}} \frac{1}{n} \sum_{\lambda \in m} \left(2 + \overline{\delta}_{n,p_\lambda}^{(\text{pen}W)} \right) \sigma_\lambda^2$$

avec $\overline{\delta}_{n,p_\lambda}^{(\text{pen}W)} \rightarrow 0$ lorsque $np_\lambda \rightarrow +\infty$.

⁴Précisons que la définition proposée pour $C_{W,\infty}$ n'a pas toujours du sens, dans la mesure où $R_{2,W}$ dépend parfois autant de $n\widehat{p}_\lambda$ que de n , et que l'on ne peut faire tendre $n\widehat{p}_\lambda$ vers l'infini en gardant n fixé. Pour être parfaitement rigoureux, il faut donc ne considérer $C_{W,\infty}$ définie que pour certains exemples de poids, par les valeurs données dans la Table 2.

Concentration. En utilisant la formule close (36), on peut montrer une inégalité de concentration pour la pénalité par rééchantillonnage [4, Proposition 3], dont nous énonçons ici une version simplifiée :

Proposition 7. *Soit $\gamma > 0$, $A_n \geq 2$, W un vecteur de poids échangeable indépendant de $(X_i, Y_i)_{1 \leq i \leq n}$. Soit m une partition de \mathcal{X} et $\text{pen}_W(m)$ définie par (35). On suppose que*

(Ab) *Les données sont bornées : $\|Y_i\|_\infty \leq A < \infty$.*

(An) *Le niveau de bruit est uniformément minoré : $\sigma(X_i) \geq \sigma_{\min} > 0$ p.s.*

Alors, une constante absolue $L_1 > 0$ et une constante $L_2(A/\sigma_{\min}, \gamma) > 0$ existent telles qu'avec probabilité au moins $1 - L_1 n^{-\gamma}$,

$$\begin{aligned} & |\text{pen}_W(m) - \mathbb{E}[\text{pen}_W(m) \mid \mathcal{P}_m]| \mathbb{1}_{\min_{\lambda \in m} \{n\hat{p}_\lambda\} \geq A_n} \leq L_2 C_W \\ & \times \sup_{np \geq A_n} \{R_{1,W}(n, p) + R_{2,W}(n, p)\} \frac{\ln(n)}{\sqrt{A_n D_m}} \mathbb{E}[p_2(m)] . \end{aligned} \quad (41)$$

En particulier, pour les poids Efr(n), Poi(1), Rad(1/2), Rho($n/2$) et Loo, une constante absolue $L_3 > 0$ existe telle que sur le même événement,

$$|\text{pen}_W(m) - \mathbb{E}[\text{pen}_W(m) \mid \mathcal{P}_m]| \mathbb{1}_{\min_{\lambda \in m} \{n\hat{p}_\lambda\} \geq A_n} \leq \frac{C_W}{C_{W,\infty}} \frac{L_2 L_3 \ln(n)}{\sqrt{A_n D_m}} \mathbb{E}[p_2(m)] .$$

On remarque que la borne (41) sur les déviations de pen_W est réduite d'un facteur $1/\sqrt{A_n}$ par rapport à la borne (18)–(19) sur les déviations de la pénalité idéale. Même s'il ne s'agit que de majorations, on retrouve une trace du phénomène de sur-concentration pour les estimateurs par rééchantillonnage que l'on avait mis en évidence dans un cas simple en Section 4.4.

6.2. Inégalité-oracle. L'événement $\min_{\lambda \in m} \{n\hat{p}_\lambda\} \geq A_n$ ayant une grande probabilité (avec $A_n \propto \ln(n)$) dès que $\min_{\lambda \in m} \{np_\lambda\} \geq L \ln(n)$, la Proposition 7, combinée à la concentration de pen_{id} (Proposition 3) et aux calculs d'espérance (14) et (40), conduit à une inégalité-oracle via le Lemme 2 du premier cours. Plus précisément, on obtient le théorème suivant [4, Théorème 1].

Théorème 3. *Soit \mathcal{M}_n une famille de partitions de \mathcal{X} et $(X_i, Y_i)_{1 \leq i \leq n}$ un échantillon i.i.d. satisfaisant les hypothèses suivantes :*

(P1) *Famille de complexité polynômiale : $\text{Card}(\mathcal{M}_n) \leq c_{\mathcal{M}} n^{\alpha_{\mathcal{M}}}$.*

(P2) *Richesse minimale de \mathcal{M}_n : $\exists m_0 \in \mathcal{M}_n$ t.q. $D_{m_0} \in [\sqrt{n}; c_{\text{rich}} \sqrt{n}]$.*

(Ab) *Les données sont bornées : $\|Y_i\|_\infty \leq A < \infty$.*

(An) *Le niveau de bruit est uniformément minoré : $\sigma(X_i) \geq \sigma_{\min} > 0$ p.s.*

(Ap) *Décroissance polynômiale du biais : il existe $\beta_1 \geq \beta_2 > 0$ et $C_b^+, C_b^- > 0$ telles que*

$$\forall m \in \mathcal{M}_n , \quad C_b^- D_m^{-\beta_1} \leq \ell(s^*, s_m^*) \leq C_b^+ D_m^{-\beta_2} .$$

(Ar $_{\ell}^X$) Régularité inférieure des partitions $m \in \mathcal{M}_n$ vis-à-vis de $\mathcal{L}(X)$: il existe $c_{r,\ell}^X > 0$ telle que

$$\forall m \in \mathcal{M}_n, \quad D_m \min_{\lambda \in m} p_{\lambda} \geq c_{r,\ell}^X .$$

Soit pen_W la pénalité définie par (35) avec $C_W = C_{W,\infty}$ (les poids étant choisis parmi $\text{Efr}(n)$, $\text{Poi}(1)$, $\text{Rad}(1/2)$, $\text{Rho}(n/2)$ et Loo).

Alors, il existe une constante $K_1 > 0$ telle qu'avec probabilité au moins $1 - K_1 n^{-2}$, pour tout

$$\widehat{m} \in \operatorname{argmin}_{m \in \mathcal{M}_n \text{ t.q. } \min_{\lambda \in m} \{n\widehat{p}_{\lambda}\} \geq 3} \{P_n \gamma(\widehat{s}_m) + \text{pen}_W(m)\} ,$$

on a

$$\ell(s^*, \widehat{s}_{\widehat{m}}) \leq \left(1 + (\ln(n))^{-1/5}\right) \inf_{m \in \mathcal{M}_n} \{\ell(s^*, \widehat{s}_m)\} . \quad (42)$$

De plus,

$$\mathbb{E}[\ell(s^*, \widehat{s}_{\widehat{m}})] \leq \left(1 + (\ln(n))^{-1/5}\right) \mathbb{E}\left[\inf_{m \in \mathcal{M}_n} \{\ell(s^*, \widehat{s}_m)\}\right] + \frac{A^2 K_1}{n^2} . \quad (43)$$

La constante K_1 peut dépendre des constantes apparaissant dans les hypothèses **(Ab)**, **(An)**, **(Ap)**, **(Ar $_{\ell}^X$)**, **(P1)** et **(P2)** mais pas de n .

Pour une discussion des hypothèses du Théorème 3, en particulier **(Ap)**, et des jeux d'hypothèses alternatifs, nous renvoyons à [4, Section 3.3].

6.3. Adaptation. L'inégalité-oracle fournie par le Théorème 3 est l'occasion d'illustrer comment l'on peut construire un estimateur adaptatif via une procédure de sélection de modèles vérifiant une inégalité-oracle.

On suppose $\mathcal{X} = [0, 1]^d$ pour simplifier, et l'on considère la famille $\mathcal{M}_n^{(\text{reg})}$ des partitions régulières de \mathcal{X} avec un nombre de morceaux maximal $M_n = n$ (voir Section 1.3). Dans la lignée des hypothèses du Théorème 3, on suppose également :

(Ab) Les données sont bornées : $\|Y_i\|_{\infty} \leq A < \infty$.

(An) Le niveau de bruit est uniformément minoré : $\sigma(X_i) \geq \sigma_{\min} > 0$ p.s.

(Ad $_{\ell}$) La densité de X est bornée inférieurement :

$$\exists c_X^{\min} > 0, \quad \forall I \subset \mathcal{X}, \quad \mathbb{P}(X \in I) \geq c_X^{\min} \text{Leb}(I) .$$

(Ah) La fonction de régression $s^* = \eta$ est α -hölderienne avec $\alpha \in (0; 1]$: il existe $R > 0$ tel que

$$s^* \in \mathcal{H}(\alpha, R) \quad \text{c'est-à-dire} \quad \forall x_1, x_2 \in \mathcal{X}, \quad |s^*(x_1) - s^*(x_2)| \leq R \|x_1 - x_2\|_{\infty}^{\alpha} .$$

On définit l'estimateur

$$\widetilde{s} := \widehat{s}_{\widehat{m}} \quad \text{avec} \quad \widehat{m} \in \operatorname{argmin}_{m \in \mathcal{M}_n^{(\text{reg})} \text{ t.q. } \min_{\lambda \in m} \{n\widehat{p}_{\lambda}\} \geq 3} \{P_n \gamma(\widehat{s}_m) + \text{pen}_W(m)\}$$

et pen_W la pénalité définie par (35) avec $C_W = C_{W,\infty}$ (les poids étant choisis parmi $\text{Efr}(n)$, $\text{Poi}(1)$, $\text{Rad}(1/2)$, $\text{Rho}(n/2)$ et Loo). Alors, d'après le Théorème 3, pour majorer le risque de \tilde{s} , il suffit de majorer le risque de l'oracle

$$\begin{aligned} & \inf_{m \in \mathcal{M}_n^{(\text{reg})}} \{ \mathbb{E} [\ell (s^*, \widehat{s}_m)] \} \\ &= \inf_{1 \leq T \leq n^{1/d}} \left\{ \ell (s^*, s_{m_T}^*) + \frac{1}{n} \sum_{\lambda \in m_T} (1 + \delta_{n,p_\lambda}) \sigma_\lambda^2 \right\} \\ &\leq \inf_{1 \leq T \leq n^{1/d} (c_X^{\min})^{1/d}} \left\{ \|s^* - s_{m_T}^*\|_\infty^2 + \frac{D_{m_T}}{n} (1 + \kappa_2) \left(\sigma_{\max}^2 + \|s^* - s_{m_T}^*\|_\infty^2 \right) \right\} \end{aligned}$$

Or, $D_{m_T} = T^d$ et comme $s^* \in \mathcal{H}(\alpha, R)$, on a

$$\|s^* - s_{m_T}^*\|_\infty \leq RT^{-\alpha}$$

et donc

$$\begin{aligned} & \inf_{m \in \mathcal{M}_n^{(\text{reg})}} \{ \mathbb{E} [\ell (s^*, \widehat{s}_m)] \} \\ &\leq \inf_{1 \leq T \leq n^{1/d} (c_X^{\min})^{1/d}} \left\{ R^2 T^{-2\alpha} + \frac{T^d}{n} (1 + \kappa_2) (\sigma_{\max}^2 + R^2 T^{-2\alpha}) \right\} \\ &\leq (2 + \kappa_2) \inf_{1 \leq T \leq n^{1/d} (c_X^{\min})^{1/d}} \left\{ R^2 T^{-2\alpha} + \frac{\sigma_{\max}^2 T^d}{n} \right\} \\ &\leq 2(2 + \kappa_2) R^{\frac{2d}{2\alpha+d}} n^{\frac{-2\alpha}{2\alpha+d}} \sigma_{\max}^{\frac{4\alpha}{2\alpha+d}} \end{aligned}$$

en choisissant

$$T = T_0 \approx \left(\frac{R^2 n}{\sigma_{\max}} \right)^{\frac{1}{2\alpha+d}}.$$

On en déduit que

$$\mathbb{E} [\ell (s^*, \tilde{s})] \leq K_2 R^{\frac{2d}{2\alpha+d}} n^{\frac{-2\alpha}{2\alpha+d}} \sigma_{\max}^{\frac{4\alpha}{2\alpha+d}} + \frac{K_3 A^2}{n^2} \quad (44)$$

pour des constantes $K_2 > 0$ numérique et $K_3(d, A, c_X^{\min}, R, \alpha, \sigma_{\min}) > 0$.

Nous renvoyons à [4, Théorème 2] pour un énoncé et une preuve précis de ce résultat, où il est également prouvé que si le niveau de bruit est régulier :

(**A** σ) σ est K_σ -Lipschitz par morceaux, avec au plus J_σ sauts,

alors, on peut supprimer l'hypothèse (**An**) et l'on a

$$\mathbb{E} [\ell (s^*, \tilde{s})] \leq K_2 R^{\frac{2d}{2\alpha+d}} n^{\frac{-2\alpha}{2\alpha+d}} \|\sigma\|_{L^2(\text{Leb})}^{\frac{4\alpha}{2\alpha+d}} + \frac{K_4 A^2}{n^2} \quad (45)$$

$$\text{avec } \|\sigma\|_{L^2(\text{Leb})} := \left[\frac{1}{\text{Leb}(\mathcal{X})} \int_{\mathcal{X}} \sigma^2(t) dt \right]^{1/2} > 0$$

et $K_4 = K_4(d, A, c_X^{\min}, R, \alpha, K_\sigma, J_\sigma) > 0$.

Les majorations du risque (44) et (45) montrent que l'estimateur \tilde{s} est adaptatif, au sens où :

- dans le cas homoscédastique, lorsque $s^* \in \mathcal{H}(\alpha, R)$ avec $\alpha > 0$ (avec une définition appropriée lorsque $\alpha > 1$) et $\mathcal{X} \subset \mathbb{R}^d$, la vitesse minimax d'estimation de s^* pour $\|\cdot\|_{L^2}^2$ est

$$R^{\frac{2d}{2\alpha+d}} n^{\frac{-2\alpha}{2\alpha+d}} \sigma^{\frac{4\alpha}{2\alpha+d}} ,$$

à un facteur multiplicatif près ne dépendant pas de n , R et σ [49, 34, 52]. Ainsi, aucun estimateur ne peut atteindre une vitesse plus rapide que (44) (à une constante près) uniformément sur toutes les lois P telles que $s^* \in \mathcal{H}(\alpha, R)$ avec $\alpha \in]0, 1]$.

- dans le cas hétéroscédastique, la dépendance du risque minimax en $\sigma(\cdot)$ n'est connue que sous des hypothèses plus fortes. Lorsque $d = 1$ et $\alpha = 1$ et $\sigma(\cdot)$ est suffisamment régulière, la majoration (45) correspond aux minorations minimax à une constante multiplicative près indépendante de $\sigma(\cdot)$ [18, 26].
- l'estimateur \tilde{s} est construit sans utiliser la connaissance exacte de $R, \alpha, \sigma(\cdot)$, condition nécessaire pour pouvoir parler d'adaptation.

Notons que l'exposant de n n'est pas suffisant pour donner du sens à une majoration telle que (44) en pratique : en effet, à horizon n fixé, on ne peut jamais distinguer une fonction $s^* \in \mathcal{C}^\infty$ d'une fonction très irrégulière, la valeur de la constante R pour laquelle $s^* \in \mathcal{H}(\alpha, R)$ est tout aussi importante. La majoration (44) peut en effet être plus fine en choisissant α plus petit que le niveau de régularité de s^* si cela permet de réduire la constante R .

Insistons pour finir sur le fait que les pénalités par rééchantillonnage n'ont pas été construites dans le but spécifique de s'adapter à l'hétéroscélasticité. On peut notamment citer à ce sujet une discussion d'un article de Wu [51] par Efron :

«The jackknife and bootstrap are general-purpose devices, not specifically adapted to take advantage of a special model like

$$y = X\beta + e , \quad \text{var}(e) = \Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2) .$$

Comparisons with specially adapted methods (...) are misleading if this is not made clear.»

En particulier, ici, il n'y a pas lieu de chercher à comparer l'estimateur \tilde{s} à des estimateurs *ad hoc*, construits avec cet objectif particulier [18, 26, 28, 29], et qui ont toutes les chances de se comporter un peu mieux que \tilde{s} pour l'analyse de données dont on sait qu'elles sont hétéroscédastiques. L'étude ci-dessus a essentiellement pour but d'illustrer la capacité naturelle du rééchantillonnage à s'adapter à des propriétés inconnues des données telles que l'hétéroscélasticité, y compris dans un cadre de sélection de modèles.

Il est donc assez tentant de conjecturer que les pénalités par rééchantillonnage s'adaptent effectivement à de nombreuses propriétés inconnues de la loi P des données pour de nombreux problèmes de sélection d'estimateurs. Les résultats mentionnés en Section 7 plaident également en ce sens.

6.4. **Comparaison des poids.** Nous renvoyons à [4, Sections 4 et 6.2] pour la comparaison des différents poids de rééchantillonnage.

7. ESTIMATION DE DENSITÉ PAR MOINDRES CARRÉS

Les résultats présentés pour l'estimation de densité sont extraits de [40] pour le cas indépendant et de [36] pour le cas dépendant. Nous renvoyons à ces deux références pour plus d'informations, ainsi qu'à [35] dont l'introduction présente ces résultats de manière synthétique.

RÉFÉRENCES

- [1] David J. Aldous. Exchangeability and related topics. In *École d'été de probabilités de Saint-Flour, XIII—1983*, volume 1117 of *Lecture Notes in Math.*, pages 1–198. Springer, Berlin, 1985.
- [2] Sylvain Arlot. *Resampling and Model Selection*. PhD thesis, University Paris-Sud 11, December 2007. <http://tel.archives-ouvertes.fr/tel-00198803/>.
- [3] Sylvain Arlot. *V-fold cross-validation improved : V-fold penalization*, February 2008. arXiv :0802.0566v2.
- [4] Sylvain Arlot. Model selection by resampling penalization. *Electron. J. Stat.*, 3 :557–624 (electronic), 2009.
- [5] Sylvain Arlot. Choosing a penalty for model selection in heteroscedastic regression, June 2010. arXiv :0812.3141v2.
- [6] Sylvain Arlot and Pascal Massart. Data-driven calibration of penalties for least-squares regression. *J. Mach. Learn. Res.*, 10 :245–279 (electronic), 2009.
- [7] Philippe Barbe and Patrice Bertail. *The Weighted Bootstrap*, volume 98 of *Lecture Notes in Statistics*. Springer-Verlag, New York, 1995.
- [8] Peter L. Bartlett, Stéphane Boucheron, and Gábor Lugosi. Model selection and error estimation. *Machine Learning*, 48 :85–113, 2002.
- [9] Peter L. Bartlett, Olivier Bousquet, and Shahar Mendelson. Local Rademacher complexities. *Ann. Statist.*, 33(4) :1497–1537, 2005.
- [10] Rudolf Beran. The impact of the bootstrap on statistical algorithms and theory. *Statist. Sci.*, 18(2) :175–184, 2003. Silver anniversary of the bootstrap.
- [11] Dennis D. Boos. Introduction to the bootstrap world. *Statist. Sci.*, 18(2) :168–174, 2003. Silver anniversary of the bootstrap.
- [12] Stéphane Boucheron, Olivier Bousquet, and Gábor Lugosi. Theory of classification : a survey of some recent advances. *ESAIM Probab. Stat.*, 9 :323–375 (electronic), 2005.
- [13] Leo Breiman. The little bootstrap and other methods for dimensionality selection in regression : X -fixed prediction error. *J. Amer. Statist. Assoc.*, 87(419) :738–754, 1992.
- [14] Leo Breiman. Random forests. *Machine Learning*, 45 :5–32, 2001. 10.1023/A :1010933404324.
- [15] Peter Bühlmann and Bin Yu. Analyzing bagging. *Ann. Statist.*, 30(4) :927–961, 2002.
- [16] George Casella, editor. *Silver Anniversary of the Bootstrap*. Institute of Mathematical Statistics, Bethesda, MD, 2003. *Statist. Sci.* **18** (2003), no. 2.
- [17] Thomas J. DiCiccio and Bradley Efron. Bootstrap confidence intervals. *Statist. Sci.*, 11(3) :189–228, 1996. With comments and a rejoinder by the authors.

- [18] Sam Efromovich and Mark Pinsker. Sharp-optimal and adaptive estimation for heteroscedastic nonparametric regression. *Statist. Sinica*, 6(4) :925–942, 1996.
- [19] Bradley Efron. Bootstrap methods : another look at the jackknife. *Ann. Statist.*, 7(1) :1–26, 1979.
- [20] Bradley Efron. *The Jackknife, the Bootstrap and Other Resampling Plans*, volume 38 of *CBMS-NSF Regional Conference Series in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, Pa., 1982.
- [21] Bradley Efron. Estimating the error rate of a prediction rule : improvement on cross-validation. *J. Amer. Statist. Assoc.*, 78(382) :316–331, 1983.
- [22] Bradley Efron. Second thoughts on the bootstrap. *Statist. Sci.*, 18(2) :135–140, 2003. Silver anniversary of the bootstrap.
- [23] Bradley Efron and Robert J. Tibshirani. *An Introduction to the Bootstrap*, volume 57 of *Monographs on Statistics and Applied Probability*. Chapman and Hall, New York, 1993.
- [24] Bradley Efron and Robert J. Tibshirani. Improvements on cross-validation : the .632+ bootstrap method. *J. Amer. Statist. Assoc.*, 92(438) :548–560, 1997.
- [25] Magalie Fromont. Model selection by bootstrap penalization for classification. *Mach. Learn.*, 66(2–3) :165–207, 2007.
- [26] Leonid Galtchouk and Sergey Pergamenshchikov. Adaptive asymptotically efficient estimation in heteroscedastic nonparametric regression via model selection, October 2008. arXiv :0810.1173v1.
- [27] Yongchao Ge, Sandrine Dudoit, and Terence P. Speed. Resampling-based multiple testing for microarray data analysis. *Test*, 12(1) :1–77, 2003. With comments and a rejoinder by the authors.
- [28] Xavier Gendre. Simultaneous estimation of the mean and the variance in heteroscedastic Gaussian regression. *Electron. J. Stat.*, 2 :1345–1372, 2008.
- [29] Xavier Gendre. *Model selection in heteroscedastic regression*. PhD thesis, University Nice Sophia-Antipolis, June 2009. <http://tel.archives-ouvertes.fr/tel-00397608/>.
- [30] Peter Hall. Asymptotic properties of the bootstrap for heavy-tailed distributions. *Ann. Probab.*, 18(3) :1342–1360, 1990.
- [31] Peter Hall. *The Bootstrap and Edgeworth Expansion*. Springer Series in Statistics. Springer-Verlag, New York, 1992.
- [32] Vladimir Koltchinskii. Local Rademacher complexities and oracle inequalities in risk minimization. *Ann. Statist.*, 34(6) :2593–2656, 2006.
- [33] Vladimir Koltchinskii and Dmitriy Panchenko. Rademacher processes and bounding the risk of function learning. In *High dimensional probability, II (Seattle, WA, 1999)*, volume 47 of *Progr. Probab.*, pages 443–457. Birkhäuser Boston, Boston, MA, 2000.
- [34] A. P. Korostelëv and A. B. Tsybakov. *Minimax Theory of Image Reconstruction*, volume 82 of *Lecture Notes in Statistics*. Springer-Verlag, New York, 1993.
- [35] Matthieu Lerasle. *Rééchantillonnage et sélection de modèles optimale pour l'estimation de la densité de variables indépendantes ou mélangeantes*. PhD thesis, INSA de Toulouse, June 2009.
- [36] Matthieu Lerasle. Optimal model selection for stationary data under various mixing conditions. arXiv :0911.1497, October 2010.
- [37] Enno Mammen. *When Does Bootstrap Work ? Asymptotic Results and Simulations*, volume 77 of *Lecture Notes in Statistics*. Springer, 1992.
- [38] David M. Mason and Michael A. Newton. A rank statistics approach to the consistency of a general bootstrap. *Ann. Statist.*, 20(3) :1611–1624, 1992.

- [39] Pascal Massart. *Concentration Inequalities and Model Selection*, volume 1896 of *Lecture Notes in Mathematics*. Springer, Berlin, 2007. Lectures from the 33rd Summer School on Probability Theory held in Saint-Flour, July 6–23, 2003, With a foreword by Jean Picard.
- [40] Matthieu Lerasle. Optimal model selection in density estimation. arXiv :0910.1654v2, 2009.
- [41] Annette M. Molinaro, Richard Simon, and Ruth M. Pfeiffer. Prediction error estimation : a comparison of resampling methods. *Bioinformatics*, 21(15) :3301–3307, 2005.
- [42] Dimitris N. Politis. The impact of bootstrap methods on time series analysis. *Statist. Sci.*, 18(2) :219–230, 2003. Silver anniversary of the bootstrap.
- [43] Dimitris N. Politis, Joseph P. Romano, and Michael Wolf. *Subsampling*. Springer Series in Statistics. Springer-Verlag, New York, 1999.
- [44] Jens Præstgaard and Jon A. Wellner. Exchangeably weighted bootstraps of the general empirical process. *Ann. Probab.*, 21(4) :2053–2086, 1993.
- [45] Jun Shao. Bootstrap model selection. *J. Amer. Statist. Assoc.*, 91(434) :655–665, 1996.
- [46] Jun Shao and Dong Sheng Tu. *The Jackknife and Bootstrap*. Springer Series in Statistics. Springer-Verlag, New York, 1995.
- [47] Ritei Shibata. Bootstrap estimate of Kullback-Leibler information for model selection. *Statist. Sinica*, 7(2) :375–394, 1997.
- [48] Kesar Singh. On the asymptotic accuracy of Efron’s bootstrap. *Ann. Statist.*, 9(6) :1187–1195, 1981.
- [49] Charles J. Stone. Optimal rates of convergence for nonparametric estimators. *Ann. Statist.*, 8(6) :1348–1360, 1980.
- [50] Aad W. van der Vaart and Jon A. Wellner. *Weak Convergence and Empirical Processes*. Springer Series in Statistics. Springer-Verlag, New York, 1996. With applications to statistics.
- [51] Chien-Fu Jeff Wu. Jackknife, bootstrap and other resampling methods in regression analysis. *Ann. Statist.*, 14(4) :1261–1350, 1986. With discussion and a rejoinder by the author.
- [52] Yuhong Yang and Andrew Barron. Information-theoretic determination of minimax rates of convergence. *Ann. Statist.*, 27(5) :1564–1599, 1999.
- [53] George Alastair Young. Bootstrap : more than a stab in the dark? *Statist. Sci.*, 9(3) :382–415, 1994. With discussion and a rejoinder by the author.

CNRS – ÉQUIPE SIERRA, LABORATOIRE D’INFORMATIQUE DE L’ÉCOLE NORMALE SUPÉRIEURE, (CNRS/ENS/INRIA UMR 8548), INRIA - 23 AVENUE D’ITALIE - CS 81321, 75214 PARIS CEDEX 13 - FRANCE

E-mail address: `sylvain.arlotRETIRERCECI@ens.fr`

URL: `http://www.di.ens.fr/~arlot/`