

**Sélection de modèles et sélection d'estimateurs pour
l'Apprentissage statistique (Cours Peccot)**

**Premier cours: Apprentissage statistique et sélection
d'estimateurs**

SYLVAIN ARLOT
CNRS – ÉQUIPE SIERRA

TABLE DES MATIÈRES

1. Le problème de l'apprentissage statistique	2
1.1. Cadre général	2
1.2. Exemple : prédiction	2
1.3. Exemple : régression	2
1.4. Exemple alternatif : régression sur un plan d'expérience fixe	3
1.5. Autres exemples	4
2. Estimateurs	4
2.1. Définition générale	4
2.2. Consistance, No Free Lunch	5
2.3. Exemples : Estimateurs par minimum de contraste	5
2.4. Exemple : Estimateurs des moindres carrés	5
2.5. Autres exemples d'estimateurs	5
3. Sélection d'estimateurs	6
3.1. Sélection d'estimateurs pour l'estimation	6
3.2. Sélection d'estimateurs pour l'identification	7
3.3. Décomposition biais-variance du risque : choix de modèles général	8
3.4. Décomposition biais-variance du risque : régression homoscédastique sur un design fixe, moindres carrés	9
3.5. Principe d'une sélection à l'aide des données	9
3.6. Pénalité idéale : régression homoscédastique sur un design fixe, moindres carrés	10
3.7. Autres procédures fondées sur l'estimation sans biais du risque	11
4. Une inégalité-oracle pour la sélection de modèles	11
4.1. Début de preuve d'inégalité-oracle	11
4.2. Une inégalité oracle pour la régression Gaussienne homoscédastique	11
5. Références mentionnées à la fin de l'exposé	14

1. LE PROBLÈME DE L'APPRENTISSAGE STATISTIQUE

On utilise ici le cadre général décrit dans [3].

1.1. Cadre général. On observe des variables aléatoires $\xi_1, \dots, \xi_n \in \Xi$ de loi commune P (les observations), supposées indépendantes. L'objectif est d'estimer à l'aide de l'échantillon $D_n = (\xi_i)_{1 \leq i \leq n}$ une quantité-cible s^* qui dépend de la distribution inconnue P , par exemple, la densité de P relativement à une mesure de référence μ , ou la fonction de régression associée à P .

Soit \mathbb{S} l'ensemble des valeurs possibles pour s^* . La qualité d'un candidat $t \in \mathbb{S}$ à estimer s^* est mesurée par sa *perte* $\mathcal{L}(t)$, où $\mathcal{L} : \mathbb{S} \mapsto \mathbb{R}$ est appelée *fonction de perte* (loss function). On suppose que la fonction de perte est minimale en $t = s^*$.

Plusieurs fonctions de perte peuvent être considérées pour un problème statistique donné. Dans ce cours, on se limitera au cas important où la fonction de perte s'écrit

$$\mathcal{L}(t) = \mathcal{L}_P(t) = P\gamma(t) := \mathbb{E}_{\xi \sim P} [\gamma(t; \xi)] \quad , \quad (1)$$

où $\gamma : \mathbb{S} \times \Xi \mapsto [0, \infty[$ est appelée *contraste*. Pour tout $t \in \mathbb{S}$, $P\gamma(t)$ mesure l'"incohérence" moyenne entre t et une nouvelle observation ξ de loi P .

Étant donnée la fonction de perte, on définit la *perte relative*

$$\ell(s^*, t) := \mathcal{L}_P(t) - \mathcal{L}_P(s^*) \geq 0 \quad .$$

1.2. Exemple : prédiction. Le problème de la prédiction a pour objectif de «prédire» une quantité d'intérêt $Y \in \mathcal{Y}$ à partir d'une variable explicative $X \in \mathcal{X}$ et d'un échantillon $(X_1, Y_1), \dots, (X_n, Y_n)$. Autrement dit, $\Xi = \mathcal{X} \times \mathcal{Y}$, \mathbb{S} est l'ensemble des applications mesurables $\mathcal{X} \mapsto \mathcal{Y}$ et le contraste $\gamma(t; (x, y))$ mesure une distance entre y et la valeur prédite $t(x)$. Deux cas particuliers classiques de ce problème sont la régression et la classification, détaillés ci-dessous.

1.3. Exemple : régression. On parle de régression lorsque \mathcal{Y} est un ensemble continu, c'est-à-dire $\mathcal{Y} \subset \mathbb{R}$ (ou \mathbb{R}^k en régression multivariée), avec le plus souvent $\mathcal{X} \subset \mathbb{R}^\ell$. Soit η la fonction de régression, c'est-à-dire, $\eta(x) = \mathbb{E}_{(X,Y) \sim P} [Y | X = x]$, de telle sorte que

$$\forall i \in \{1, \dots, n\} \quad , \quad Y_i = \eta(X_i) + \varepsilon_i \quad \text{avec} \quad \mathbb{E}[\varepsilon_i | X_i] = 0 \quad .$$

Un contraste souvent utilisé en régression est le *contraste des moindres carrés* $\gamma(t; (x, y)) = (t(x) - y)^2$, qui atteint son minimum sur \mathbb{S} en $t = s^* = \eta$, et la perte relative vaut

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

On peut remarquer que l'excès de risque de t est la carré de la distance $L^2(P)$ entre t et s^* , si bien que *prédiction* et *estimation* sont ici des objectifs équivalents.

Preuve de (2). Soit $t \in \mathbb{S}$, (X, Y) une v.a. de loi P et $\varepsilon = Y - \eta(X)$. Alors, par définition de η , $\mathbb{E}[\varepsilon | X] = 0$ p.s., et donc

$$P\gamma(t) = \mathbb{E} \left[(\varepsilon + \eta(X) - t(X))^2 \right] = \mathbb{E}[\varepsilon^2] + \mathbb{E} \left[(\eta(X) - t(X))^2 \right] + 2\mathbb{E}[\varepsilon(\eta(X) - t(X))]$$

et le dernier terme est nul car, P -presque sûrement,

$$\mathbb{E}[\varepsilon(\eta(X) - t(X)) | X] = (\eta(X) - t(X)) \mathbb{E}[\varepsilon | X] = 0 .$$

La perte $P\gamma(t)$ est donc minimale pour $t = \eta = s^*$, et la perte relative a bien l'expression donnée par (2). \square

1.4. Exemple alternatif : régression sur un plan d'expérience fixe.

On considère également souvent le cadre de la régression en supposant que le plan d'expérience (en anglais, *design*) X_1, \dots, X_n est déterministe. On le note alors souvent en lettres minuscules x_1, \dots, x_n par convention.

L'échantillon se résume alors à $Y = (Y_1, \dots, Y_n) \in \mathbb{R}^n$ avec

$$\forall i \in \{1, \dots, n\}, \quad Y_i = \eta(x_i) + \varepsilon_i$$

avec $\varepsilon_1, \dots, \varepsilon_n$ des v.a. indépendantes, centrées. Lorsque les ε_i sont i.i.d., on parle de régression *homoscédastique*; le cas contraire est appelé *hétéroscédastique*.

Si les Y_i ne sont plus de même loi, on peut tout de même considérer le risque et l'excès de risque correspondant à la loi $P = \mathcal{L}(X_{n+1}, Y_{n+1})$ définie par $X_{n+1} \sim \mathcal{U}(x_1, \dots, x_n)$ et $\mathcal{L}(Y_{n+1} | X_{n+1} = x_i) = \mathcal{L}(Y_i)$. En notant $F = (F_i)_{1 \leq i \leq n} = (\eta(x_i))_{1 \leq i \leq n}$ et $\varepsilon = (\varepsilon_i)_{1 \leq i \leq n}$, on peut alors résumer le problème à

$$Y = F + \varepsilon \in \mathbb{R}^n$$

L'équivalent du contraste des moindres carrés est alors de chercher $t \in \mathbb{S} = \mathbb{R}^n$ tel que la perte quadratique

$$P\gamma(t) = \mathbb{E}_Y \left[\frac{1}{n} \|t - Y\|^2 \right] = \mathbb{E}_Y \left[\frac{1}{n} \sum_{i=1}^n (t_i - Y_i)^2 \right]$$

est minimale. Or,

$$\begin{aligned} \mathbb{E}_Y \left[\frac{1}{n} \|t - Y\|^2 \right] &= \mathbb{E}_Y \left[\frac{1}{n} \|t - F - \varepsilon\|^2 \right] \\ &= \frac{1}{n} \mathbb{E}_Y \left[\|t - F\|^2 + \|\varepsilon\|^2 - 2 \langle \varepsilon, t - F \rangle \right] \\ &= \frac{1}{n} \mathbb{E}_Y \left[\|t - F\|^2 + \|\varepsilon\|^2 \right] \end{aligned}$$

TAB. 1. Régression : design fixe vs. design aléatoire.

	Design aléatoire	Design fixe
Données	$(X_i, Y_i)_{1 \leq i \leq n}$ i.i.d. $\sim P$	$Y = F + \varepsilon \in \mathbb{R}^n$
$\mathcal{D}(X_{n+1}, Y_{n+1})$	$\frac{P}{P}$	$\mathcal{D}(x_J, Y_J)$ avec $J \sim \mathcal{U}(1, \dots, n)$
$t \in \mathbb{S}$	$t : \mathcal{X} \rightarrow \mathbb{R}$	$t \in \mathbb{R}^n$
Perte de t	$P\gamma(t) = \mathbb{E}_{(X,Y) \sim P} \left[(Y - t(X))^2 \right]$	$\mathbb{E}_Y \left[\frac{1}{n} \ Y - t\ ^2 \right] = \mathbb{E}_Y \left[\frac{1}{n} \sum_{i=1}^n (Y_i - t_i)^2 \right]$
Cible s^*	$\eta : x \rightarrow \mathbb{E}[Y X = x]$	$F = (\eta(x_1), \dots, \eta(x_n))$
Perte relative $\ell(s^*, t)$	$\mathbb{E}_{(X,Y) \sim P} \left[(t(X) - \eta(X))^2 \right]$	$\frac{1}{n} \ F - t\ ^2 = \frac{1}{n} \sum_{i=1}^n (\eta(x_i) - t_i)^2$
Risque empirique de t	$P_n\gamma(t) = \frac{1}{n} \sum_{i=1}^n (Y_i - t(X_i))^2$	$\ Y - t\ ^2 = \sum_{i=1}^n (Y_i - t_i)^2$

si bien que la perte quadratique est minimale pour $t = F = s^*$ et la perte relative vaut

$$\ell(s^*, t) = \frac{1}{n} \|t - F\|^2 = \frac{1}{n} \sum_{i=1}^n (t_i - F_i)^2 .$$

Le Tableau 1 résume les correspondances entre quantités-clés entre les cadres de régression à désign déterministe et aléatoire.

La différence entre considérer les X_i aléatoires ou déterministes est souvent assez faible (même si les philosophies sous-jacentes et les objectifs sont assez différents), sauf si l'on pose des hypothèses trop fortes sur le design (e.g., ne considérer que $x_i = i/n$, qui est un cas trop éloigné de X_i i.i.d.). Voir par exemple [5, 6] où l'on montre un résultat à design fixe pour en déduire un à design aléatoire.

Du point de vue de l'analyse théorique, on choisit celui des deux qui simplifie le plus l'analyse. Ainsi, nous nous focaliserons surtout sur le cas du design fixe dans ce cours et dans le suivant, tandis que nous supposerons essentiellement le design aléatoire lors des deux derniers cours (ce qui permet de traiter plus simplement l'hétéroscédasticité et les méthodes de rééchantillonnage).

1.5. Autres exemples. Pour une présentation formelle d'autres exemples (estimation de densité, classification), nous renvoyons à la Section 1.2 de [3]. Voir aussi [12] pour un survol des résultats récents en classification.

2. ESTIMATEURS

2.1. Définition générale. On définit un *algorithme statistique* (ou règle d'apprentissage) \mathcal{A} comme étant une application mesurable $\mathcal{A} : \bigcup_{n \in \mathbb{N}} \Xi^n \mapsto \mathbb{S}$. Étant donné un échantillon $D_n = (\xi_i)_{1 \leq i \leq n} \in \Xi^n$, la sortie de \mathcal{A} , $\mathcal{A}(D_n) = \hat{s}^{\mathcal{A}}(D_n) \in \mathbb{S}$, est un estimateur de s^* . La qualité de \mathcal{A} est alors mesurée par sa perte $\mathcal{L}_P(\hat{s}^{\mathcal{A}}(D_n))$, que l'on souhaite minimiser.

On définit alors le *risque* de \mathcal{A} pour un échantillon de taille n et de loi P par

$$\mathbb{E}_{\xi \sim P^{\otimes n}} \left[P\gamma(\hat{s}^{\mathcal{A}}(D_n)) \right] = R(\mathcal{A}, P, n) .$$

Du point de vue statistique, on parle plutôt de l'estimateur $\widehat{s} = \widehat{s}(D_n) = \widehat{s}^{\mathcal{A}}(D_n) = \mathcal{A}(D_n) \in \mathbb{S}$, en oubliant sa relation fonctionnelle \mathcal{A} avec l'échantillon D_n .

L'objectif est de construire $\widehat{s}(D_n) \in S$ tel que $\ell(s^*, \widehat{s}(D_n))$ est petit en espérance, ou avec grande probabilité.

2.2. Consistance, No Free Lunch. Nous renvoyons à [13] à propos des différentes notions de consistance et des résultats du type «No Free Lunch Theorems». Voir aussi [1].

2.3. Exemples : Estimateurs par minimum de contraste. La famille des *estimateurs par minimum de contraste* est classique en statistique. Étant donné un modèle $S \subset \mathbb{S}$, on dit que $\widehat{s}(D_n)$ est un estimateur par minimum de contraste sur S s'il minimise p.s. sur S le contraste empirique

$$t \mapsto P_n \gamma(t) = \frac{1}{n} \sum_{i=1}^n \gamma(t; \xi_i) \quad \text{où} \quad P_n = \frac{1}{n} \sum_{i=1}^n \delta_{\xi_i} .$$

L'heuristique sous-jacente est que l'espérance du contraste empirique $P_n \gamma(t)$ est $P \gamma(t)$, qui est minimale pour $t = s^*$. On peut donc espérer que minimiser $P_n \gamma(t)$ sur $S \subset \mathbb{S}$ fournit un bon estimateur de s^* , pourvu que S soit bien choisi.

Outre les estimateurs des moindres carrés (voire la sous-section suivante), un exemple classique est l'estimateur du maximum de vraisemblance en estimation de densité, qui correspond au contraste «log-vraisemblance» $\gamma(t; x) = -\ln(t(x))$.

Pour le problème de la prédiction, un estimateur minimisant le contraste empirique $P_n \gamma(t)$ est plus couramment appelé *minimiseur du risque empirique* [21].

2.4. Exemple : Estimateurs des moindres carrés. On prend $\gamma(t; (x, y)) = (t(x) - y)^2$ en régression (on pourrait aussi prendre le contraste des moindres carrés en estimation de densité), et le modèle S est un sous-espace vectoriel de \mathbb{S} . Par exemple, si S est l'espace des fonctions constantes par morceaux associé à une partition fixée de \mathcal{X} , alors on obtient un *régressogramme*. On peut aussi construire un modèle comme le sous-espace vectoriel généré par les premiers éléments d'une base telle que la base de Fourier ou une base d'ondelettes.

2.5. Autres exemples d'estimateurs. Pour les autres exemples d'estimateurs présentés à l'oral, nous renvoyons aux références suivantes :

- pour la régression ridge (à noyau) et les splines de lissage : [22, 18, 2],
- pour les SVM : [18, 11],
- pour le Lasso et autres méthodes de régularisation L^1 : [19, 25, 4],
- pour les k -plus proches voisins : [13].

3. SÉLECTION D'ESTIMATEURS

Un grand nombre d'estimateurs peuvent être utilisés pour résoudre un problème statistique donné. Soit $(\mathcal{A}_m)_{m \in \mathcal{M}}$ une famille d'algorithmes statistiques et $(\widehat{s}_m(D_n))_{m \in \mathcal{M}}$ les estimateurs correspondants. Le problème de *sélection d'estimateur* (ou sélection d'algorithme statistique) est celui de choisir, à l'aide des données uniquement, l'un de ces algorithmes, c'est-à-dire choisir $\widehat{m}(D_n) \in \mathcal{M}$. L'estimateur final de s^* est alors $\widehat{s}_{\widehat{m}(D_n)}(D_n)$. La difficulté principale est alors que l'on ne dispose pas de données indépendantes de D_n pour évaluer la perte de chacun de $\widehat{s}_m(D_n)$: il faut utiliser D_n à nouveau.

Un exemple fondamental de ce problème est la sélection de modèles, où l'on dispose d'une famille $(S_m)_{m \in \mathcal{M}_n}$ de modèles, et pour tout $m \in \mathcal{M}_n$, \widehat{s}_m est un estimateur par minimum de contraste qui lui est associé (par exemple, l'estimateur des moindres carrés en régression). Le problème de la sélection de modèles a fait l'objet de nombreux travaux, voir notamment [16].

Tout d'abord, précisons l'objectif visé par une procédure de choix d'estimateur : estimation ou identification (voir aussi à ce sujet les sections 2.2 à 2.4 de [3]).

3.1. Sélection d'estimateurs pour l'estimation. On parle d'objectif d'*estimation* (ou de prédiction) lorsque l'on souhaite minimiser la perte relative $\ell(s^*, \widehat{s}_{\widehat{m}}(D_n))$ de l'estimateur final. Le choix optimal est alors appelé *oracle*, défini par

$$m^* := m^*(D_n) \in \arg \min_{m \in \mathcal{M}_n} \{ \ell(s^*, \widehat{s}_m(D_n)) \} . \quad (3)$$

L'oracle $m^*(D_n)$ dépendant de la distribution inconnue P des données, on ne peut espérer en général sélectionner exactement l'oracle avec grande probabilité. On se fixe donc un objectif plus atteignable, qui est de faire presque aussi bien que l'oracle en termes de perte relative. Selon le cadre, on peut le formaliser de différentes manières.

D'une part, dans un cadre asymptotique, on dit qu'une procédure de choix d'estimateur $\widehat{m}(\cdot)$ est *efficace* ou *asymptotiquement optimale* lorsque

$$\frac{\ell(s^*, \widehat{s}_{\widehat{m}(D_n)}(D_n))}{\inf_{m \in \mathcal{M}_n} \{ \ell(s^*, \widehat{s}_m(D_n)) \}} \xrightarrow[n \rightarrow \infty]{p.s.} 1 .$$

D'autre part, dans un cadre non-asymptotique, on dit qu'une procédure de choix d'estimateur $\widehat{m}(\cdot)$ satisfait une *inégalité-oracle* avec constante C_n (et terme de reste R_n) lorsque

$$\ell(s^*, \widehat{s}_{\widehat{m}(D_n)}(D_n)) \leq C_n \inf_{m \in \mathcal{M}_n} \{ \ell(s^*, \widehat{s}_m(D_n)) \} + R_n \quad (4)$$

est vérifié en espérance, ou avec grande probabilité (c'est-à-dire, une probabilité au moins aussi grande que $1 - C'/n^2$, pour une constante $C' > 0$). Remarquons que si (4) est vérifiée pour n assez grand sur un événement de probabilité $1 - C'/n^2$, et si C_n tend vers 1 quand n tend vers l'infini, alors \widehat{m} est asymptotiquement optimale.

L'oracle est parfois défini comme $\arg \min_{m \in \mathcal{M}_n} \{ \mathbb{E} [\ell (s^*, \widehat{s}_m(D_n))] \}$, conduisant à une forme affaiblie de l'inégalité-oracle (4) :

$$\mathbb{E} [\ell (s^*, \widehat{s}_{\widehat{m}(D_n)}(D_n))] \leq C_n \inf_{m \in \mathcal{M}_n} \{ \mathbb{E} [\ell (s^*, \widehat{s}_m(D_n))] \} + R_n . \quad (5)$$

En effet, si (4) a lieu en espérance, alors (5) a lieu car

$$\mathbb{E} \left[\inf_{m \in \mathcal{M}_n} \{ \ell (s^*, \widehat{s}_m(D_n)) \} \right] \leq \inf_{m \in \mathcal{M}_n} \{ \mathbb{E} [\ell (s^*, \widehat{s}_m(D_n))] \} .$$

Notons enfin que les procédures de sélection d'estimateurs pour l'estimation sont souvent utilisées pour construire des *estimateurs adaptatifs* (par exemple au sens du minimax) en choisissant une famille d'estimateurs adéquate [7].

3.2. Sélection d'estimateurs pour l'identification. On parle d'objectif d'*identification* lorsque l'on souhaite identifier la «meilleure» règle d'apprentissage (c'est-à-dire, celle qui minimise le risque) avec une probabilité maximale. C'est un objectif plus ambitieux que l'objectif d'estimation, dans la mesure où sélectionner l'oracle avec grande probabilité implique une inégalité-oracle (5) avec la même probabilité.

En retour, parler d'identification nécessite de faire une hypothèse supplémentaire, qui est que cette «meilleure règle» soit *identifiable*, c'est-à-dire bien définie de manière unique.

Considérons le cas de la sélection de modèles pour préciser les choses (voir par exemple [24] pour une description plus formelle du cas général).

L'objectif est alors de sélectionner le «vrai modèle» S_{m_0} , défini comme le plus petit modèle (au sens de l'inclusion) parmi $(S_m)_{m \in \mathcal{M}_n}$ qui contient s^* . En particulier, il faut supposer que $s^* \in \bigcup_{m \in \mathcal{M}_n} S_m$ (pour n assez grand), et que m_0 est effectivement identifiable (c'est-à-dire qu'il existe un m_0 tel que tout modèle S_m contenant s^* contient nécessairement S_{m_0}).

On mesure alors la qualité d'une procédure de choix de modèles pour l'identification par la probabilité de sélectionner m_0 (que l'on souhaite maximiser). L'équivalent de l'optimalité asymptotique est appelée *consistance* (souvent appelée «consistance en modèle» pour ne pas la confondre avec la consistance au sens de la perte), définie par

$$\mathbb{P} (\widehat{m}(D_n) = m_0) \xrightarrow[n \rightarrow \infty]{} 1 .$$

Comme nous l'avons remarqué, lorsqu'il existe un unique «vrai modèle» m_0 , une procédure de choix de modèles consistante (en modèle) est nécessairement asymptotiquement optimale. On peut naturellement se demander s'il est possible de construire une procédure de choix de modèles atteignant simultanément l'objectif d'identification (lorsque $s^* \in \bigcup_{m \in \mathcal{M}_n} S_m$) et l'objectif d'estimation (lorsque $s^* \notin \bigcup_{m \in \mathcal{M}_n} S_m$).

Yang [23] a prouvé que c'est impossible en régression en général, en montrant qu'aucune procédure de choix de modèles ne peut être simultanément

adaptative au sens du minimax (une propriété qui découle de l'optimalité asymptotique lorsque la famille de modèles est bien choisie) et consistante en modèle.

En revanche, en ajoutant des hypothèses sur le problème considéré, il est parfois possible de posséder les deux propriétés simultanément (voir [20] et l'introduction de [23]).

3.3. Décomposition biais-variance du risque : choix de modèles général. À partir de cette section et jusqu'à la fin de ce premier cours, nous nous focaliserons sur le problème de la sélection de modèles, étant entendu que les principales idées développées ici pourraient être étendues au problème plus général de la sélection d'estimateurs. Nous reparlerons du problème général au cours des chapitres suivants.

Soit $(S_m)_{m \in \mathcal{M}_n}$ une famille de modèles, i.e., pour tout $m \in \mathcal{M}_n$, $S_m \subset \mathbb{S}$, et notons \hat{s}_m l'estimateur par minimum de contraste associé (γ étant un contraste fixé).

On peut décomposer l'excès de risque $\mathbb{E}[\ell(s^*, \hat{s}_m(D_n))]$ en deux termes. D'une part, comme $\hat{s}_m(D_n) \in S_m$ p.s., l'excès de risque est minoré par l'*erreur d'approximation* (aussi appelée *biais*) définie par

$$\ell(s^*, S_m) := \inf_{t \in S_m} \{\ell(s^*, t)\} .$$

Lorsque cet infimum est atteint en au moins un élément de S_m , on note $\ell(s^*, s_m^*) = \ell(s^*, S_m)$. Nous ferons souvent cette hypothèse simplificatrice dans la suite, sachant que l'on peut également définir s_m^* comme un δ -minimiseur de $\ell(s^*, t)$ sur S_m avec δ choisi assez petit (e.g., $\delta = n^{-2}$). Le biais quantifie la «distance» entre le modèle S_m et la cible s^* . Il est d'autant plus petit que S_m est grand.

Le second terme, appelée *erreur d'estimation* (ou *variance*), est défini de manière générale par

$$\mathbb{E}[\ell(s^*, \hat{s}_m(D_n))] - \ell(s^*, S_m) = \mathbb{E}[P\gamma(\hat{s}_m(D_n))] - \inf_{t \in S_m} P\gamma(t) .$$

Lorsque s_m^* est bien défini, on peut le réécrire

$$\mathbb{E}[P\gamma(\hat{s}_m(D_n)) - P\gamma(s_m^*)] .$$

Ce terme quantifie la difficulté d'estimation au sein du modèle S_m . Il est en général d'autant plus grand que le modèle S_m est grand. Dans plusieurs cadres classiques, on peut prouver qu'il est (approximativement) proportionnel au nombre de paramètres libres du modèle S_m .

Au final, on a la décomposition de l'excès de risque

$$\begin{aligned} \mathbb{E}[\ell(s^*, \hat{s}_m(D_n))] &= \ell(s^*, S_m) + \mathbb{E}[P\gamma(\hat{s}_m(D_n)) - P\gamma(s_m^*)] \\ &= \text{Biais} + \text{Variance} . \end{aligned} \quad (6)$$

Un choix de modèle optimal en terme d'excès de risque nécessite donc de trouver un compromis entre ces deux termes, i.e., choisir S_m assez grand pour diminuer son biais, mais pas trop pour éviter le sur-apprentissage (lorsque le terme de variance devient prépondérant). On parle souvent de *compromis biais-variance*. Typiquement, pour l'oracle, les termes de biais et de variance sont très proches.

3.4. Décomposition biais-variance du risque : régression homoscedastique sur un design fixe, moindres carrés. Dans ce cadre simple, on peut calculer simplement le risque. Rappelons que l'on a $Y = F + \varepsilon$ et que l'on considère l'estimateur $\widehat{F}_m = A_m Y$ où A_m est la matrice de projection orthogonale sur l'espace vectoriel S_m . Alors,

$$\frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 = \frac{1}{n} \|(A_m - I)F\|^2 + \frac{1}{n} \|A_m \varepsilon\|^2 \quad (7)$$

$$\begin{aligned} \text{d'où } \mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right] &= \frac{1}{n} \|(A_m - I)F\|^2 + \frac{\sigma^2 \operatorname{tr}(A_m^\top A_m)}{n} \\ &= \frac{1}{n} \|(A_m - I)F\|^2 + \frac{\sigma^2 \dim(S_m)}{n} \\ &= \text{Erreur d'approximation} + \text{Erreur d'estimation} \end{aligned}$$

si les ε_i sont indépendants, centrés et de variance σ^2 , en utilisant que $A_m^\top A_m = A_m$ car A_m est une matrice de projection orthogonale, et que $\operatorname{tr}(A_m) = \dim(S_m)$. Dans cette preuve, on a utilisé le Lemme suivant.

Lemme 1. Soient $M \in \mathcal{M}_n(\mathbb{R})$ et $\varepsilon_1, \dots, \varepsilon_n$ des v.a. centrées et de même variance σ^2 . Alors,

$$\mathbb{E}[\langle \varepsilon, M\varepsilon \rangle] = \sigma^2 \operatorname{tr}(M) .$$

Démonstration.

$$\begin{aligned} \mathbb{E}[\langle \varepsilon, M\varepsilon \rangle] &= \mathbb{E} \left[\sum_{i=1}^n \varepsilon_i (M\varepsilon)_i \right] = \mathbb{E} \left[\sum_{i=1}^n \varepsilon_i \sum_{j=1}^n M_{i,j} \varepsilon_j \right] \\ &= \sum_{i=1}^n \sum_{j=1}^n (M_{i,j} \mathbb{E}[\varepsilon_i \varepsilon_j]) = \sum_{i=1}^n (M_{i,i} \mathbb{E}[\varepsilon_i^2]) = \sigma^2 \operatorname{tr}(M) . \end{aligned}$$

□

3.5. Principe d'une sélection à l'aide des données. Pour choisir $\widehat{m}(D_n) \in \mathcal{M}_n$, on utilise la stratégie générale suivante :

$$\widehat{m}(D_n) \in \arg \min_{m \in \mathcal{M}_n} \{ \operatorname{crit}(m; D_n) \} ,$$

où $\operatorname{crit}(m; D_n)$ est un critère pouvant dépendre des données, idéalement égal au risque $P\gamma(\widehat{s}_m(D_n))$ (à constante additive près).

3.5.1. Principe de l'estimation sans biais du risque. On prend

$$\operatorname{crit}(m; D_n) \approx \mathbb{E}[P\gamma(\widehat{s}_m(D_n))]$$

(à constante additive près). Cela fonctionne s'il n'y a pas trop de modèles (ou d'estimateurs), modulo une inégalité de concentration de $P\gamma(\widehat{s}_m(D_n))$ autour de son espérance (uniforme en $m \in \mathcal{M}$). Si \mathcal{M}_n est trop riche, on ne peut pas obtenir cette concentration uniforme avec une simple borne d'union, si bien qu'il peut être nécessaire de faire différemment (i.e., utiliser la structure de la famille d'estimateurs, ou bien regrouper astucieusement les m par groupes).

3.5.2. Pénalisation. Le risque empirique est optimiste comme estimation du risque (exemple d'une famille de modèles emboîtés, sur-apprentissage).

Pour estimer sans biais le risque, il faut donc lui ajouter une quantité appelée pénalité. Une procédure de choix de modèle par pénalisation est une procédure de la forme

$$\widehat{m}(D_n) \in \arg \min_{m \in \mathcal{M}_n} \{ P_n \gamma(\widehat{s}_m(D_n)) + \text{pen}(m; D_n) \} . \quad (8)$$

Le choix idéal (à constante additive près) pour la pénalité est appelé *pénalité idéale*, définie par

$$\text{pen}_{\text{id}}(m; D_n) := (P - P_n) \gamma(\widehat{s}_m(D_n)) .$$

En appliquant le principe d'estimation sans biais du risque, on aboutit à la pénalité déterministe idéale

$$\mathbb{E} [\text{pen}_{\text{id}}(m; D_n)] = \mathbb{E} [(P - P_n) \gamma(\widehat{s}_m(D_n))] .$$

3.6. Pénalité idéale : régression homoscedastique sur un design fixe, moindres carrés. Reprenons l'exemple de la régression homoscedastique sur un design fixe, avec la perte quadratique. Le risque empirique vaut

$$\frac{1}{n} \left\| \widehat{F}_m - Y \right\|^2 = \frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 - \frac{2}{n} \langle (A_m - I_n)F, \varepsilon \rangle - \frac{2}{n} \langle A_m \varepsilon, \varepsilon \rangle + \frac{1}{n} \|\varepsilon\|^2 ,$$

$$\text{d'où } \mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - Y \right\|^2 \right] = \mathbb{E} \left[\frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 \right] - \frac{2\sigma^2 \text{tr}(A_m)}{n} + \sigma^2$$

si les ε_i sont indépendants, centrés et de variance σ^2 .

On en déduit l'expression suivante pour la pénalité idéale et son espérance (en lui ajoutant $n^{-1} \|\varepsilon\|^2$ qui ne dépend pas de m) :

$$\begin{aligned} \text{pen}_{\text{id}}(m) &= \frac{1}{n} \left\| \widehat{F}_m - F \right\|^2 - \frac{1}{n} \left\| \widehat{F}_m - Y \right\|^2 + \frac{1}{n} \|\varepsilon\|^2 \\ &= \frac{2}{n} \langle A_m \varepsilon, \varepsilon \rangle + \frac{2}{n} \langle (A_m - I_n)F, \varepsilon \rangle \\ \mathbb{E} [\text{pen}_{\text{id}}(m)] &= \frac{2\sigma^2 \text{tr}(A_m)}{n} = \frac{2\sigma^2 \dim(S_m)}{n} \end{aligned} \quad (9)$$

La pénalité $\frac{2\sigma^2 \dim(S_m)}{n}$ est aussi appelée C_p de Mallows. Si l'on n'avait pas utilisé les propriétés de A_m (matrice de projection orthogonale), on aurait aboutit à la pénalité plus générale $\frac{2\sigma^2 \text{tr}(A_m)}{n}$, aussi appelée C_L de Mallows [15]. Pour les estimateurs linéaires, $\text{tr}(A_m)$ est appelé nombre de degrés de

liberté généralisés associé à A_m , par analogie avec le cas des estimateurs par projection orthogonale.

3.7. Autres procédures fondées sur l'estimation sans biais du risque.

Outre les exemples que nous verrons aux cours suivants, nous renvoyons à la Section 3 de [3] pour des références concernant d'autres procédures de choix de modèles (ou d'estimateurs) fondées sur l'estimation sans biais du risque.

4. UNE INÉGALITÉ-ORACLE POUR LA SÉLECTION DE MODÈLES

Le but de cette section est de prouver une inégalité-oracle à l'aide d'un schéma de preuve très général, qui sera utilisé dans les cours suivants, et qui est notamment pour l'essentiel celui utilisé par [7, 16].

4.1. Début de preuve d'inégalité-oracle.

Lemme 2. Soit $\text{pen} : \mathcal{M}_n \mapsto \mathbb{R}$ une pénalité (pouvant dépendre des données). Soit $(A(m))_{m \in \mathcal{M}_n}$, $(B(m))_{m \in \mathcal{M}_n} \in \mathbb{R}^{\mathcal{M}_n}$. Sur l'événement Ω où pour tout $m, m' \in \mathcal{M}_n$,

$$(\text{pen}(m) - \text{pen}_{\text{id}}(m, D_n)) - (\text{pen}(m') - \text{pen}_{\text{id}}(m', D_n)) \leq A(m) + B(m') , \quad (10)$$

$$\text{on a } \forall \hat{m} \in \arg \min_{m \in \mathcal{M}_n} \{P_n \gamma(\hat{s}_m(D_n)) + \text{pen}(m)\} , \quad (11)$$

$$\ell(s^*, \hat{s}_{\hat{m}}(D_n)) - B(\hat{m}) \leq \inf_{m \in \mathcal{M}_n} \{\ell(s^*, \hat{s}_m(D_n)) + A(m)\} . \quad (12)$$

Preuve du Lemme 2. D'après la définition (11) de \hat{m} , pour tout $m \in \mathcal{M}_n$,

$$P_n \gamma(\hat{s}_{\hat{m}}(D_n)) + \text{pen}(\hat{m}, D_n) \leq P_n \gamma(\hat{s}_m(D_n)) + \text{pen}(m)$$

ce qui peut se réécrire : $\forall m \in \mathcal{M}_n$,

$$\ell(s^*, \hat{s}_{\hat{m}}(D_n)) + \text{pen}(\hat{m}) - \text{pen}_{\text{id}}(\hat{m}, D_n) \leq \ell(s^*, \hat{s}_m(D_n)) + \text{pen}(m) - \text{pen}_{\text{id}}(m, D_n) . \quad (13)$$

En utilisant l'équation (10) avec $m' = \hat{m}$, (13) entraîne

$$\forall m \in \mathcal{M}_n, \quad \ell(s^*, \hat{s}_{\hat{m}}(D_n)) \leq \ell(s^*, \hat{s}_m(D_n)) + A(m) + B(\hat{m}) ,$$

d'où le résultat. \square

4.2. Une inégalité oracle pour la régression Gaussienne homoscedastique. En s'appuyant sur les calculs précédents, le Lemme 2 et deux inégalités de concentration, on peut reprover un résultat dû à Birgé et Massart [10, 9].

Théorème 1. On se place dans le cadre de la régression à design fixe, avec la perte des moindres carrés. On suppose donnée une famille d'estimateurs des moindres carrés $(\hat{F}_m)_{m \in \mathcal{M}_n}$ associée à une famille de modèles $(S_m)_{m \in \mathcal{M}_n}$ qui sont des s.e.v. de dimension finie de \mathbb{R}^n . On suppose que $\text{Card}(\mathcal{M}_n)$ est fini.

On suppose que les données sont de la forme $Y = F + \varepsilon$ avec $F \in \mathbb{R}^n$ inconnu et $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ (cadre Gaussien homoscedastique). On note A_m

la matrice de projection orthogonale sur S_m , de telle sorte que $\widehat{F}_m = A_m Y$ et $\text{tr}(A_m) = \dim(S_m) = D_m$.

Soit $K > 1$ une constante. Pour tout $m \in \mathcal{M}_n$, on définit

$$\text{pen}(m) = \frac{K\sigma^2 \dim(S_m)}{n} .$$

Alors, il existe $C(K) > 0$ telle que pour tout $x \geq 0$, avec probabilité au moins $1 - 4 \text{Card}(\mathcal{M}_n)e^{-x}$, on a pour tout

$$\widehat{m} \in \arg \min_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \|Y - \widehat{F}_m\|^2 + \text{pen}(m) \right\}$$

l'inégalité-oracle : pour tout $\delta > 0$,

$$\frac{1}{n} \|\widehat{F}_{\widehat{m}} - F\|^2 \leq \left(\frac{1 + (K-2)_+}{1 - (2-K)_+} + \delta \right) \inf_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \|\widehat{F}_m - F\|^2 \right\} + \frac{C(K)x\sigma^2}{\delta n} \quad (14)$$

Notons que l'on a d'après la preuve du Théorème 1 la majoration

$$C(K) \leq L \max \{ (K-1)^{-3}, (K-1)^2 \} ,$$

avec $L > 0$ une constante numérique.

En particulier, en prenant $K = 2$ (pénalité C_p de Mallows), (14) s'écrit

$$\frac{1}{n} \|\widehat{F}_{\widehat{m}} - F\|^2 \leq (1 + \delta) \inf_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \|\widehat{F}_m - F\|^2 \right\} + \frac{Lx\sigma^2}{\delta n} , \quad (15)$$

qui implique l'optimalité asymptotique de \widehat{m} (en intégrant par rapport à x ou en prenant $x = 2 \ln(n) + \ln(\text{Card}(\mathcal{M}_n))$), et en faisant tendre δ vers 0 à une vitesse adéquate et en supposant que la vitesse d'apprentissage de l'oracle est non-paramétrique :

$$\inf_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \|\widehat{F}_m - F\|^2 \right\} \gg \frac{1}{n} .$$

Pour prouver le Théorème 1, on a besoin des deux inégalités de concentration suivantes [16].

Proposition 3. Soit ξ un vecteur Gaussien standard de \mathbb{R}^n , $\alpha \in \mathbb{R}^n$, $M \in \mathcal{M}_n(\mathbb{R})$. On note $\|\alpha\|$ la norme euclidienne de α , $\|M\|$ la norme d'opérateur de M (sup de modules des valeurs propres). Alors, pour tout $x \geq 0$,

$$\begin{aligned} \mathbb{P} \left(|\langle \xi, \alpha \rangle| \leq \sqrt{2x} \|\alpha\|_2 \right) &\geq 1 - 2e^{-x} \\ \mathbb{P} \left(|\langle \xi, M\xi \rangle - \text{tr}(M)| \leq 2\sqrt{x \text{tr}(M^\top M)} + 2\|M\| x \right) &\geq 1 - 2e^{-x} . \end{aligned}$$

Preuve du Théorème 1. On découpe cette preuve en 5 étapes principales :

- (1) Pour chaque $x \geq 0$ et $m \in \mathcal{M}_n$, on concentre $\langle A_m \varepsilon, \varepsilon \rangle$ autour de $\sigma^2 \text{tr}(A_m)$ et $\langle (A_m - I_n)F, \varepsilon \rangle$ autour de zéro. Ceci définit un événement $\Omega_{m,x}$ de probabilité $1 - 4e^{-x}$, sur lequel pour tout $\theta > 0$

$$|\langle A_m \varepsilon, \varepsilon \rangle - \sigma^2 D_m| \leq \theta D_m \sigma^2 + (\theta^{-1} + 2) x \sigma^2 \quad (16)$$

$$|\langle (A_m - I_n)F, \varepsilon \rangle| \leq \theta \|(A_m - I_n)F\|^2 + \frac{1}{2\theta} x\sigma^2 \quad (17)$$

en utilisant le fait que pour tout $a, b > 0$, $2\sqrt{ab} \leq \theta a + \theta^{-1}b$.

- (2) Soit $\Omega_x = \bigcap_{m \in \mathcal{M}} \Omega_{m,x}$, pour lequel une borne d'union donne $\mathbb{P}(\Omega) \geq 1 - 4 \text{Card}(\mathcal{M}_n)e^{-x}$. Sur Ω_x , en utilisant (16), (17) et (9), pour tout $m \in \mathcal{M}_n$ et $\theta > 0$,

$$\begin{aligned} \left| \text{pen}_{\text{id}}(m) - \frac{2\sigma^2 D_m}{n} \right| &\leq \frac{2\theta\sigma^2 D_m}{n} + \frac{2\theta}{n} \|(A_m - I_n)F\|^2 + (4 + 3\theta^{-1}) \frac{x\sigma^2}{n} \\ &= 2\theta \mathbb{E} \left[\frac{\|F - \widehat{F}_m\|^2}{n} \right] + (4 + 3\theta^{-1}) \frac{x\sigma^2}{n} . \end{aligned} \quad (18)$$

- (3) On majore ce terme de reste à l'aide de la perte relative $\|F - \widehat{F}_m\|^2$ directement (et non son espérance). En effet, en utilisant (7) et le fait que A_m est une matrice de projection orthogonale, on a

$$\begin{aligned} \|F - \widehat{F}_m\|^2 &= \|(A_m - I)F\|^2 + \|A_m \varepsilon\|^2 = \|(A_m - I)F\|^2 + \langle A_m \varepsilon, \varepsilon \rangle \\ &\geq \|(A_m - I)F\|^2 + \sigma^2 D_m - \theta\sigma^2 D_m - (\theta^{-1} + 2) x\sigma^2 \\ &\geq (1 - \theta) \mathbb{E} \left[\|F - \widehat{F}_m\|^2 \right] - (\theta^{-1} + 2) x\sigma^2 \end{aligned}$$

sur Ω_x d'après (16). Donc, pour tout $\theta \in]0, 1[$,

$$\mathbb{E} \left[\|F - \widehat{F}_m\|^2 \right] \leq \frac{1}{1 - \theta} \|F - \widehat{F}_m\|^2 + \frac{1}{1 - \theta} (\theta^{-1} + 2) x\sigma^2 . \quad (19)$$

En injectant (19) dans (18), on obtient ainsi pour tout $m \in \mathcal{M}_n$ et $\theta \in]0, 1[$,

$$\left| \text{pen}_{\text{id}}(m) - \frac{2\sigma^2 D_m}{n} \right| \leq \frac{2\theta}{1 - \theta} \|F - \widehat{F}_m\|^2 + \left(\frac{2\theta}{1 - \theta} (\theta^{-1} + 2) + 4 + 3\theta^{-1} \right) x\sigma^2 .$$

- (4) D'après (5) et (19), l'hypothèse du Lemme 2 est donc vérifiée sur Ω_x avec

$$\begin{aligned} A(m) &= \frac{(K - 2)_+ + 2\theta}{1 - \theta} \|F - \widehat{F}_m\|^2 + C_1(\theta, K) \frac{x\sigma^2}{n} \\ B(m) &= \frac{(2 - K)_+ + 2\theta}{1 - \theta} \|F - \widehat{F}_m\|^2 + C_2(\theta, K) \frac{x\sigma^2}{n} \end{aligned}$$

$$\text{avec } C_1(\theta, K) = 4 + 3\theta^{-1} + \frac{2\theta + (K - 2)_+}{1 - \theta} (\theta^{-1} + 2)$$

$$\text{et } C_2(\theta, K) = 4 + 3\theta^{-1} + \frac{2\theta + (2 - K)_+}{1 - \theta} (\theta^{-1} + 2) .$$

Par le Lemme 2, on obtient donc que sur Ω_x , pour tout $m \in \mathcal{M}_n$ et $\theta \in]0, 1[$,

$$\left(1 - \frac{(2 - K)_+ + 2\theta}{1 - \theta} \right) \frac{1}{n} \|F - \widehat{F}_m\|^2$$

$$\leq \left(1 + \frac{(2-K)_+ + 2\theta}{1-\theta}\right) \inf_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \|F - \widehat{F}_m\|^2 \right\} + (C_1(\theta, K) + C_2(\theta, K)) \frac{x\sigma^2}{n}$$

(5) Pour en déduire le résultat annoncé, on suppose tout d'abord que $\theta \in]0, 1[$ vérifie $\theta < (K-1)/3$, si bien que $\frac{(2-K)_+ + 2\theta}{1-\theta} < 1$. On obtient alors que

$$\frac{1}{n} \|F - \widehat{F}_{\widehat{m}}\|^2 \leq C_3(K, \theta) \inf_{m \in \mathcal{M}_n} \left\{ \frac{1}{n} \|F - \widehat{F}_m\|^2 \right\} + C_4(K, \theta) \frac{x\sigma^2}{n}$$

$$\text{avec } C_3(K, \theta) = \left(1 - \frac{(2-K)_+ + 2\theta}{1-\theta}\right)^{-1} \left(1 + \frac{(2-K)_+ + 2\theta}{1-\theta}\right)$$

$$\text{et } C_4(K, \theta) = \left(1 - \frac{(2-K)_+ + 2\theta}{1-\theta}\right)^{-1} (C_1(\theta, K) + C_2(\theta, K)) .$$

Le résultat s'en déduit en choisissant un θ adéquat.

Plus précisément, en supposant que $0 < \theta < \min\{L_1, L_2(K-1)\}$ pour des constantes numériques $L_1, L_2 > 0$, des constantes numériques $L_3, L_4, L_5 > 0$ existent telles que si $K \in]1, 2]$, alors

$$C_3(K, \theta) \leq \frac{1}{K-1} \left(1 + \frac{L_3\theta}{K-1}\right)$$

$$C_4(K, \theta) \leq \frac{L_4}{(K-1)\theta}$$

et si $K \geq 2$,

$$C_3(K, \theta) \leq (K-1)(1 + L_5\theta)$$

$$C_4(K, \theta) \leq \frac{L_4(K-1)}{\theta} .$$

Il suffit donc de prendre

$$\theta = L_6\delta \min \left\{ (K-1)^2, \frac{1}{K-1} \right\}$$

pour obtenir le résultat annoncé ($L_6 > 0$ étant une constante numérique), si bien que l'on a

$$C'(K, \delta) \leq C_4(K, \theta) \leq \frac{L_7}{\delta} \max \{ (K-1)^{-3}, (K-1)^2 \} ,$$

$L_7 > 0$ étant une constante numérique.

□

5. RÉFÉRENCES MENTIONNÉES À LA FIN DE L'EXPOSÉ

- à propos de programmation dynamique pour la détection de ruptures : [8, 17]
- à propos de la pénalité idéale pour le Lasso : [14, 26]

RÉFÉRENCES

- [1] Sylvain Arlot. Classification supervisée : des algorithmes et leur calibration automatique, 2009. Notes d'un cours de troisième année à l'École Centrale Paris. <http://www.di.ens.fr/~arlot/enseign/2009Centrale/cours-classif.pdf.gz>.
- [2] Sylvain Arlot and Francis Bach. Data-driven calibration of linear estimators with minimal penalties. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 46–54, 2009.
- [3] Sylvain Arlot and Alain Celisse. A survey of cross-validation procedures for model selection. *Statist. Surv.*, 4 :40–79, 2010.
- [4] Francis Bach. Tutorial on sparse methods for machine learning (theory and algorithms), 2010. Transparents d'un cours donné à ECML. <http://www.di.ens.fr/~fbach/>.
- [5] Yannick Baraud. Model selection for regression on a fixed design. *Probab. Theory Related Fields*, 117(4) :467–493, 2000.
- [6] Yannick Baraud. Model selection for regression on a random design. *ESAIM Probab. Statist.*, 6 :127–146 (electronic), 2002.
- [7] Andrew Barron, Lucien Birgé, and Pascal Massart. Risk bounds for model selection via penalization. *Probab. Theory Related Fields*, 113(3) :301–413, 1999.
- [8] Richard E. Bellman and Stuart E. Dreyfus. *Applied dynamic programming*. Princeton University Press, Princeton, N.J., 1962.
- [9] Lucien Birgé and Pascal Massart. Minimal penalties for Gaussian model selection. *Probab. Theory Related Fields*, 138(1-2) :33–73, 2007.
- [10] Lucien Birgé and Pascal Massart. A generalized cp criterion for gaussian model selection. Technical report, Universités de Paris 6 et Paris 7, 2010. Prépublication 647, 39 pages.
- [11] Gilles Blanchard, Olivier Bousquet, and Pascal Massart. Statistical performance of support vector machines. *Ann. Statist.*, 36(2) :489–531, 2008.
- [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] Luc Devroye, László Györfi, and Gábor Lugosi. *A probabilistic theory of pattern recognition*, volume 31 of *Applications of Mathematics (New York)*. Springer-Verlag, New York, 1996.
- [14] Bradley Efron, Trevor Hastie, Iain Johnstone, and Robert Tibshirani. Least angle regression. *Ann. Statist.*, 32(2) :407–499, 2004. With discussion, and a rejoinder by the authors.
- [15] Colin L. Mallows. Some comments on C_p . *Technometrics*, 15 :661–675, 1973.
- [16] 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.
- [17] Guillem Rigai. Pruned dynamic programming for optimal multiple change-point detection. arXiv :1004.0887, April 2010.
- [18] Bernhard Scholkopf and Alexander J. Smola. *Learning with Kernels : Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, Cambridge, MA, USA, 2001.
- [19] Robert Tibshirani. Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B*, 58(1) :267–288, 1996.

- [20] T. van Erven, P. D. Grünwald, and S. de Rooij. Catching Up Faster by Switching Sooner : A Prequential Solution to the AIC-BIC Dilemma, July 2008. arXiv :0807.1005.
- [21] Vladimir Vapnik. *Estimation of Dependences Based on Empirical Data*. Springer Series in Statistics. Springer-Verlag, New York, 1982. Translated from the Russian by Samuel Kotz.
- [22] Grace Wahba. *Spline Models for Observational Data*, volume 59 of *CBMS-NSF Regional Conference Series in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1990.
- [23] Yuhong Yang. Can the strengths of AIC and BIC be shared? A conflict between model identification and regression estimation. *Biometrika*, 92(4) :937–950, 2005.
- [24] Yuhong Yang. Consistency of cross validation for comparing regression procedures. *Ann. Statist.*, 35(6) :2450–2473, 2007.
- [25] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of The Royal Statistical Society Series B*, 68(1) :49–67, 2006.
- [26] Hui Zou, Trevor Hastie, and Robert Tibshirani. On the “degrees of freedom” of the lasso. *Ann. Statist.*, 35(5) :2173–2192, 2007.

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/`