

1 Complexité de Rademacher

Exercice 1. Inégalité de Ledoux Le but de cet exercice est de prouver le théorème suivant :

Theorème 1 (Inégalité de Ledoux). Si φ est B -Lipshitz et $\varepsilon_1^n = \{\varepsilon_i\}_{i=1}^n$ est une suite i.i.d. de variables Rad(1/2), nous avons

$$\mathbb{E} \left[\sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \varphi(\theta^\top x_i) \right] \leq B \mathbb{E} \left[\sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \theta^\top x_i \right].$$

Nous allons montrer par récurrence que pour tout $k \in \{1, \dots, n\}$, **toutes fonctions** $b : \Theta \rightarrow \mathbb{R}$, $a_i : \Theta \rightarrow \mathbb{R}$, $i \in \{1, \dots, k\}$ et toutes fonctions 1-Lipshitz $\varphi_i : \mathbb{R} \rightarrow \mathbb{R}$, $i = 1, \dots, k$,

$$\mathbb{E} \left[\sup_{\theta \in \Theta} b(\theta) + \sum_{i=1}^k \varepsilon_i \varphi_i(a_i(\theta)) \right] \leq \mathbb{E} \left[\sup_{\theta \in \Theta} b(\theta) + \sum_{i=1}^k \varepsilon_i a_i(\theta) \right] \quad (1)$$

1. Pour toutes fonctions $\varphi, \psi : \Theta \rightarrow \mathbb{R}$ et $\varepsilon \sim \text{Rad}(1/2)$, montrer que :

$$\begin{aligned} \mathbb{E} \left[\sup_{\theta \in \Theta} \{\varphi(\theta) + \varepsilon \psi(\theta)\} \right] &= \frac{1}{2} \left[\sup_{\theta, \theta' \in \Theta^2} \{\varphi(\theta) + \varphi(\theta') + \psi(\theta) - \psi(\theta')\} \right] \\ &= \frac{1}{2} \left[\sup_{\theta, \theta' \in \Theta^2} \{\varphi(\theta) + \varphi(\theta') + |\psi(\theta) - \psi(\theta')|\} \right] \end{aligned} \quad (2)$$

2. Supposons que (1) est vérifiée pour $k \in \mathbb{N}$. Montrer que

$$\begin{aligned} &\mathbb{E} \left[\sup_{\theta \in \Theta} b(\theta) + \sum_{i=1}^{k+1} \varepsilon_i \varphi_i(a_i(\theta)) \right] \\ &\leq \mathbb{E} \left[\mathbb{E} \left[\sup_{\theta, \theta' \in \Theta} \frac{b(\theta) + b(\theta')}{2} + \sum_{i=1}^k \varepsilon_i \frac{\varphi_i(a_i(\theta)) + \varphi_i(a_i(\theta'))}{2} + \frac{|a_{k+1}(\theta) - a_{k+1}(\theta')|}{2} \middle| \varepsilon_1^k \right] \right] \\ &= \mathbb{E} \left[\mathbb{E} \left[\sup_{\theta \in \Theta} b(\theta) + \sum_{i=1}^k \varepsilon_i \varphi_i(a_i(\theta)) + \varepsilon_{k+1} a_{k+1}(\theta) \middle| \varepsilon_1^k \right] \right] \\ &= \mathbb{E} \left[\sup_{\theta \in \Theta} b(\theta) + \sum_{i=1}^{k+1} \varepsilon_i a_i(\theta) \right] \end{aligned}$$

3. Conclure.

Exercice 2 (Dernière étape pour calculer la borne supérieure de la complexité de Rademacher). Nous considérons $X_1^n = \{X_i\}_{i=1}^n$ un n -échantillon de distribution \mathbb{P} . Nous supposons qu'il existe $R > 0$, tel que $\mathbb{E}[\|X\|^2] \leq R^2$. On considère $\varepsilon_1^n = \{\varepsilon_i\}_{i=1}^n$ une suite de variables aléatoires indépendantes de Rademacher $\mathcal{R}(1/2)$ indépendante de X_1^n . Montrer que :

$$\mathbb{E} \left[\sup_{\|\theta\| \leq D} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \theta^\top X_i \right] \leq \frac{RD}{\sqrt{n}}. \quad (3)$$

Exercice 3 (Inégalité de McDiarmid). Le but de cet exercice est de prouver l'inégalité de McDiarmid, énoncée ci-dessous :

Dfinition 1 (Fonction de différence bornée). $g : Z^n \rightarrow \mathbb{R}$ est une fonction de différence bornée s'il existe des constantes $\{c_i\}_{i=1}^n$ telles que, pour tout $z_1^n = \{z_i\}_{i=1}^n$ et $\tilde{z}_1^n = \{\tilde{z}_i\}_{i=1}^n$,

$$|g(z_1^n) - g(\tilde{z}_1^n)| \leq \sum_{i=1}^n c_i \mathbb{1}_{\{z_i \neq \tilde{z}_i\}}.$$

Theorème 2 (Inégalité de McDiarmid). Si g est une fonction de différence bornée et Z_i sont des variables aléatoires indépendantes alors

$$\begin{aligned} \mathbb{P}(g(Z_1, \dots, Z_n) - \mathbb{E}[g(Z_1, \dots, Z_n)] \geq \varepsilon) &\leq e^{-\frac{2\varepsilon^2}{\sum_{i=1}^n c_i^2}} \\ \mathbb{P}(\mathbb{E}[g(Z_1, \dots, Z_n)] - g(Z_1, \dots, Z_n) \geq \varepsilon) &\leq e^{-\frac{2\varepsilon^2}{\sum_{i=1}^n c_i^2}} \end{aligned}$$

Pour $i \in \{1, \dots, n\}$ et $z_1^i = (z_1, \dots, z_i) \in Z^i$, nous posons $g_i(z_1^i) = \mathbb{E}[g(z_1^i, Z_{i+1}^n)]$. Par convention, nous posons $g_0 = \mathbb{E}[g(Z_1^n)]$. Notons que $g_n(z_1^n) = g(z_1^n)$, pour tout $i \in \{1, \dots, n\}$, $\mathbb{E}[g_i(Z_1^i)] = \mathbb{E}[g(Z_1^n)]$ et pour $n \geq 1$

$$g(z_1^n) = g_n(z_1^n) = g_n(z_1^n) - g_{n-1}(z_1^{n-1}) + g_{n-1}(z_1^{n-1}) = \sum_{j=1}^n \{g_j(z_1^j) - g_{j-1}(z_1^{j-1})\},$$

avec la convention que $g_0(z_1^0) = g_0$.

1. Montrer que pour tout $\lambda \in \mathbb{R}$,

$$\mathbb{E} \left[e^{\lambda \{g_n(Z_1^n) - g_0\}} \right] = \mathbb{E} \left[\mathbb{E} \left[e^{\lambda \{g_n(Z_1^n) - g_{n-1}(Z_1^{n-1})\}} \mid Z_1^{n-1} \right] e^{\lambda \{g_{n-1}(Z_1^{n-1}) - g_0\}} \right]$$

2. Montrer que pour tout $i \in \{1, \dots, n\}$, $\inf_{z \in Z} g_i(z_1^{i-1}, z) \leq g_i(z_1^i) \leq \inf_{z \in Z} g_i(z_1^{i-1}, z) + c_i$ et que

$$g_{i-1}(Z_1^{i-1}) = \mathbb{E} \left[g_i(Z_1^i) \mid Z_1^{i-1} \right], \quad \mathbb{P} - \text{p.s.}$$

3. En déduire que

$$\mathbb{E} \left[e^{\lambda \{g_n(Z_1^n) - g_{n-1}(Z_1^{n-1})\}} \mid Z_1^{n-1} \right] \leq e^{\lambda^2 c_n^2 / 8}.$$

4. Montrer que :

$$\mathbb{E} \left[e^{\lambda \{g(Z_1, \dots, Z_n) - \mathbb{E}[g(Z_1, \dots, Z_n)]\}} \right] \leq e^{\frac{\lambda^2 \sum_{i=1}^n c_i^2}{8}}.$$

5. Prouver le théorème

Nous considérons

$$\Delta_{n,01}(\mathcal{C}) = \sup_{f \in \mathcal{C}} R(f) - \widehat{R}_{n,01}(f) = \sup_{f \in \mathcal{C}} \left(\mathbb{E}[\ell_{01}(Y, f(X))] - \frac{1}{n} \sum_{i=1}^n \ell_{01}(Y_i, f(X_i)) \right).$$

6. Montrer que

$$\mathbb{P}(\Delta_n(\mathcal{C}) - \mathbb{E}[\Delta_n(\mathcal{C})] \leq \varepsilon) \geq 1 - e^{-2n\varepsilon^2}$$

2 Convexification du risque

Exercice 4 (Lemme de Zhang, (cours)). On considère le problème de classification binaire, où $y \in \{\pm 1\}$. Nous considérons également la règle de prédiction suivante : prédire $y = 1$ si $g(x) \geq 0$, et prédire $y = -1$ sinon. La perte 0-1 de classification de $g(\cdot)$ en un point (x, y) est donnée par

$$\ell_{01}(g(x), y) = \begin{cases} 1, & \text{si } -g(x)y > 0, \\ 1, & \text{si } g(x) = 0 \text{ et } y = -1, \\ 0, & \text{sinon.} \end{cases}$$

Notre objectif est de trouver un prédicteur $g(x)$ de sorte à minimiser le risque 0-1 : $R_{01}(g) = \mathbb{E}[\ell_{01}(g(X), Y)]$. Étant donné un ensemble de données d'apprentissage $\mathcal{D}_n = \{(X_i, Y_i)\}_{i=1}^n$ i.i.d. la minimisation du risque empirique 0-1 consiste à trouver g dans une classe de fonctions \mathcal{C} qui minimise

$$\widehat{R}_{n,01}(g) = \frac{1}{n} \sum_{i=1}^n \ell_{01}(g(X_i), Y_i) \quad (4)$$

1. Pourquoi la minimisation du risque empirique $\widehat{R}_{n,01}$ est elle problématique ?

Soit $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$ une fonction convexe. On considère la minimisation dans une classe de fonctions \mathcal{C} du risque empirique suivant

$$\widehat{R}_{n,\phi}(g) = \frac{1}{n} \sum_{i=1}^n \phi(-g(X_i) Y_i),$$

qui est une approximation stochastique du risque $R_\phi(g) = \mathbb{E}[\phi(-g(X)Y)]$.

2. Montrer que

$$R_\phi(g) = \mathbb{E}[\eta(X)\phi(-g(X)) + (1 - \eta(X))\phi(g(X))]$$

où $\eta(X) = \mathbb{P}(Y = +1 | X)$.

On rappelle les exemples suivants :

Nom de la perte	Perte ϕ	Nom de la méthode
Moindres carrés	$\phi_2(v) = (1 + v)^2$	Moindres carrés.
Hinge	$\phi_H(v) = \max(1 + v, 0)$	SVM
Exponentielle	$\phi_E(v) = \exp(v)$	AdaBoost
Logistique	$\phi_L(v) = \log_2(1 + \exp(v))$	Régression Logistique

3. Quelles sont les propriétés cruciales de ces fonctions qui permettront de les optimiser ?

Montrer que pour tous les exemples ci dessus (sauf ϕ_2),

$$\phi \text{ est croissante, convexe, et } \phi(v) \geq \mathbb{1}_{\{v \geq 0\}}. \quad (\text{H1})$$

On s'intéresse dans la suite à une fonction ϕ générique satisfaisant (H1). Au regard de la question 2, nous considérons la fonction $H_\phi : [0, 1] \times \mathbb{R} \rightarrow \mathbb{R}$

$$H_\phi(\eta, p) = \eta\phi(-p) + (1 - \eta)\phi(p).$$

4. Montrer que $p \mapsto H_\phi(\eta, p)$ est convexe.

5. Montrer que pour $\eta \in]0; 1[$, $p \mapsto H_\phi(\eta, p)$ admet un minimum unique dans \mathbb{R} , pour $\phi \in \{\phi_2, \phi_E, \phi_L\}$. Qu'en est il pour ϕ_H ? Qu'en est il pour $\eta \in \{0; 1\}$?

Nous supposons dans la suite qu'il existe une fonction $p_\phi^*(\eta) : [0, 1] \rightarrow \bar{\mathbb{R}}$ défini par¹

$$p_\phi^*(\eta) = \arg \min_{p \in \mathbb{R}} H_\phi(\eta, p), \quad g_\phi^*(x) = p_\phi^* \circ \eta(x)$$

et nous notons

$$H_\phi^*(\eta) = \inf_{p \in \mathbb{R}} H_\phi(\eta, p) = H_\phi(\eta, p_\phi^*(\eta)),$$

$$\Delta H_\phi(\eta, p) = H_\phi(\eta, p) - H_\phi(\eta, p_\phi^*(\eta)) = H_\phi(\eta, p) - H_\phi^*(\eta).$$

Par symétrie, nous avons $H_\phi(\eta, p) = H_\phi(1 - \eta, -p)$. Lorsque que $p_\phi^*(\eta)$ n'est pas déterminé de manière unique, on choisit un minimiseur quelconque mais on suppose dans la suite que p_ϕ^* est choisi de telle sorte que $p_\phi^*(1 - \eta) = -p_\phi^*(\eta)$. En particulier, elle implique que $p_\phi^*(1/2) = 0$.

Le prédicteur bayésien pour la perte ϕ est le prédicteur bayésien pour la perte 01

6. Montrer que : $g_\phi^* \in \arg \min_{g \in \mathbb{R}^{\mathcal{X}}} R_\phi(g)$ et pour tout g :

$$\Delta R_\phi(g) = R_\phi(g) - R_\phi(g_\phi^*) = \mathbb{E}[\Delta H_\phi(\eta(X), g(X))]$$

7. Montrer que :

Nom de la perte	p_ϕ^*	H_ϕ^*
Moindres carrés	$p_\phi^*(\eta) = 2\eta - 1$	$H_\phi^*(\eta) = 4\eta(1 - \eta)$.
Hinge	$p_\phi^*(\eta) = \text{sign}(2\eta - 1)$	$H_\phi^*(\eta) = 1 - 2\eta - 1 $.
Logistique	$p_\phi^*(\eta) = \log_2 \frac{\eta}{1-\eta}$	$H_\phi^*(\eta) = -\eta \log_2 \eta - (1 - \eta) \log_2(1 - \eta)$.

8. Montrer que pour $\phi \in \{\phi_2, \phi_H, \phi_E, \phi_L\}$, la règle de décision $\text{signe}(g_\phi^*)$ est un classifieur de Bayes pour la perte 0 - 1.

9. Supposons que ϕ est croissante, convexe et différentiable. Montrer que $\text{signe}(g_\phi^*)$ est la règle de décision bayésienne pour la perte 0-1.

Controle de l'excès de risque pour la perte 01 par l'excès de risque pour la perte ϕ .

Supposons que $p_\phi^*(\eta) > 0$ lorsque $\eta > 1/2$ et qu'il existe $c > 0$ et $s \geq 1$ tels que pour toute $\eta \in [0, 1]$,

$$|1/2 - \eta|^s \leq c^s \Delta H_\phi(\eta, 0). \quad (\text{H2})$$

L'objectif de cet partie est d'établir que pour toute fonction mesurable $g(x)$.

$$R_{01}(g) - R_{01}^* \leq 2c \Delta R_\phi(g)^{1/s},$$

où R^* est l'erreur de Bayes optimale pour la perte 0 - 1.

1. Soit $\bar{\mathbb{R}}$ la droite réelle étendue ($\bar{\mathbb{R}} = \mathbb{R} \cup \{-\infty, +\infty\}$). Nous étendons une fonction convexe $\psi : \bar{\mathbb{R}} \rightarrow \bar{\mathbb{R}}$ en définissant $\psi(\infty) = \lim_{x \rightarrow \infty} \psi(x)$ et $\psi(-\infty) = \lim_{x \rightarrow -\infty} \psi(x)$. Elle garantit que le minimiseur optimal $p_\phi^*(\eta)$ donné ci-dessous est bien défini à $\eta = 0$ ou 1 pour certaines fonctions de perte.

10. Montrer que

$$R_{01}(g) - R_{01}^* \leq 2c \left(\mathbb{E}[\mathbb{1}_{\{(2\eta(X)-1)g(X) \leq 0\}} \Delta H_\phi(\eta(X), 0)] \right)^{1/s}.$$

11. Montrer que

$$(2\eta(x) - 1)g(x) < 0 \Rightarrow \Delta H_\phi(\eta(x), 0) \leq \Delta H_\phi(\eta(x), g(x))$$

12. Montrer que la condition (H2) est vérifiée avec $c = \frac{1}{2}$, $s = 2$. pour le critère quadratique.

13. Montrer que la condition (H2) est vérifiée avec $c = \frac{1}{2}$, $s = 1$ pour la perte hinge.

14. Montrer que la condition (H2) est vérifiée avec $c = \sqrt{\ln 2/2}$, $s = 2$ pour la perte logistique.