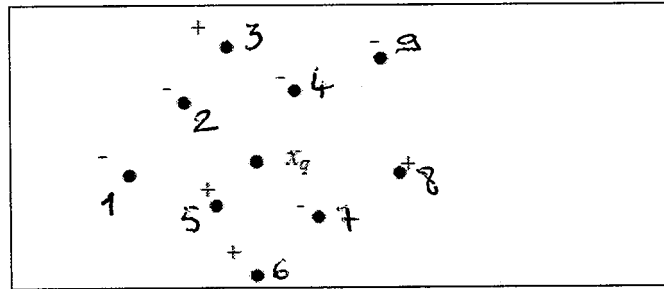


Apprentissage Statistique en Grande Dimension- Examen

Le barème est approximatif. On veillera à justifier les réponses avec soin.

Exercice 1 (4,5 pts). On considère le jeu de données dans \mathbb{R}^2 représenté dans la figure ci-dessous. On dispose de 9 points d'entraînement et d'un point x_q à prédire. Dans cet exercice, la fonction de perte



considérée est l'erreur de classification.

- (0,5 pt) Déterminez la prédiction associée au point x_q dans le cadre d'un 1-plus proche voisin et d'un 3-plus proche voisin.
- (1 pt) Calculez l'erreur de validation croisée 2-folds lorsque les points 1/2/3/4, 5/6/7/8 forment les 2 groupes et que l'on considère l'algorithme des 3-plus proches voisins.
- (1 pt) Quel est le nombre minimal de vecteurs supports que peut contenir un SVM linéaire construit sur les 9 points d'apprentissage? Expliquez votre réponse.
- (1 pt) En indiquant les abscisses et ordonnées par des x_i et y_j , fabriquer un arbre binaire générant une erreur d'entraînement nulle.
- (1 pt) Calculez l'hétérogénéité globale d'un arbre construit à partir d'une seule coupe laissant les points 2/3/4/9 ensemble, lorsque la "métrique" est l'indice de Gini.

Exercice 2 (Classification binaire, 9 pts). On suppose que (X, Y) est un couple de variables aléatoires tel que $X \in \mathcal{X}$ et $Y \in \{-1, 1\}$. On se donne une fonction $\phi : \mathbb{R} \rightarrow \mathbb{R}$ et on définit pour $f : \mathcal{X} \rightarrow \mathbb{R}$ le risque R_f par :

$$R_f = \mathbb{E}[\phi(Yf(X))].$$

On note $\eta(x) = \mathbb{P}(Y = 1|X = x)$.

- (1 pt) Exprimez R_f sous la forme $\mathbb{E}[\Psi(X)]$ où Ψ est une fonction à déterminer.
- (2 pts) On suppose que $\phi(-1) > \phi(1)$. Montrez avec soin que

$$\min_{f: \mathbb{R} \rightarrow \{-1, 1\}} R_f = \mathbb{E}[\min\{\eta(X)(\phi(1) - \phi(-1)) + \phi(-1), \eta(X)(\phi(-1) - \phi(1)) + \phi(1)\}],$$

et que ce minimum est atteint en la règle de Bayes usuelle (associée à l'erreur de classification).

- (1,5 pt) On suppose que X suit la loi uniforme sur $[0, 1]$, que $\phi(x) = \exp(-x)$ et que $\eta(x) = x$. Montrez que dans ce cas, le risque optimal satisfait :

$$\min_{f: \mathbb{R} \rightarrow \{-1, 1\}} R_f = -\frac{1}{2} \text{sh}(1) + \text{ch}(1).$$

4. A partir de maintenant, on suppose que l'on dispose d'une famille de fonctions $(f_\theta)_\theta$ de \mathcal{X} dans \mathbb{R} et que l'on cherche

$$\theta^* = \operatorname{Argmin}_{\theta \in \Theta} R_{f_\theta}$$

où θ^* est supposé exister et être unique pour simplifier.

- (a) (1 pt) Supposons que l'on dispose d'un échantillon $\{(X_1, Y_1), \dots, (X_N, Y_N)\}$ issu de (X, Y) . Définir un estimateur $\hat{\theta}_N$ de θ^* (On pensera à considérer le minimiseur d'une fonction de l'échantillon naturellement associée à R_{f_θ}).
- (b) (1,5 pt + 0,5 pt bonus) On suppose dans cette question que $\phi(x) = e^{-x}$ et on définit l'algorithme suivant :

$$\theta_{k+1} = \theta_k + \frac{\gamma}{N} \sum_{i=1}^N \nabla_{\theta} f_{\theta_k}(X_i) Y_i e^{-Y_i f_{\theta_k}(X_i)}.$$

Expliquez à quoi correspond cet algorithme et pourquoi il a vocation à converger vers $\hat{\theta}_N$ lorsque $k \rightarrow +\infty$. Cette convergence a-t-elle lieu lorsque $f_\theta(x) = \langle x, \theta \rangle$?

- (c) (1 pt) Ecrivez la version mini-batch (stochastique) de l'algorithme ci-dessus.
- (d) (1 pt) Notons Θ , l'espace des paramètres associé aux fonctions f_θ . Quels types de problèmes peut-on rencontrer en fonction de la taille Θ ? (en quelques lignes).

Exercice 3 (6,5 points). On note $\{(X_1, Y_1), \dots, (X_N, Y_N)\}$ un échantillon d'apprentissage (*i.i.d.*) de taille N tel que pour tout $1 \leq i \leq N$, $X_i \in \mathbb{R}^p$ and $Y_i \in \mathbb{R}$. On note \mathcal{L} la fonction définie par

$$\mathcal{L}(\theta) = \sum_{i=1}^N (Y_i - \langle X_i, \theta \rangle)^2 + \lambda \sum_{j=1}^{p-1} |\theta_{j+1} - \theta_j|,$$

où λ est un réel positif (Il s'agit d'un Fused-Lasso simplifié). On note (lorsque celui-ci est bien défini)

$$\hat{\theta} = \operatorname{Argmin}_{\theta \in \mathbb{R}^p} \mathcal{L}(\theta).$$

- (0,5 pt) Quel sens donnez-vous à $\hat{\theta}$? Quel type d'application cela peut-il avoir ?
- (1 pt) Montrez que le sous-gradient associé à \mathcal{L} existe puis calculez $\partial_{\theta_j} \mathcal{L}$ pour $j = 2, \dots, p-1$ (On pensera à séparer les cas).
- (1,5 pt) On suppose que $\hat{\theta}_1 < \hat{\theta}_2 < \dots < \hat{\theta}_p$ et que $\mathbf{X}^T \mathbf{X} = I_p$. Déterminez l'expression de $\hat{\theta}_j$ pour $j = 2, \dots, p-1$ dans ce cas.
- (1,5 pt) On suppose dans cette question que $|\mathbf{X}_j^T \mathbf{Y}| \leq \frac{\lambda}{N}$ pour $j = 2, \dots, p-1$, que $|\mathbf{X}_1^T \mathbf{Y}| \leq \frac{\lambda}{2N}$ et que $|\mathbf{X}_p^T \mathbf{Y}| \leq \frac{\lambda}{2N}$. Montrez qu'alors $0_{\mathbb{R}^p}$ est un minimum de \mathcal{L} (On rappelle que \mathbf{X} est la matrice ayant pour lignes X_1, \dots, X_N et colonnes $\mathbf{X}_1, \dots, \mathbf{X}_p$, et $\mathbf{Y} = (Y_1, \dots, Y_N)^T$).
- Dans cette question, on modifie légèrement la fonction \mathcal{L} en posant

$$\tilde{\mathcal{L}}(\theta) = \sum_{k=1}^N (Y_i - \langle X_i, \theta \rangle)^2 + \lambda \left(\sum_{j=1}^{p-1} |\theta_{j+1} - \theta_j| + |\theta_p| \right).$$

- (a) (0,5 pt) Montrez que $\sum_{j=1}^{p-1} |\theta_{j+1} - \theta_j| + |\theta_p| = \|D\theta\|_1$ où D est une matrice que l'on déterminera.
- (b) (1,5 pt) Supposons que $(\mathbf{X}D^{-1})^T \mathbf{X}D = I_p$. En s'appuyant sur le cours et sur un changement de variable, déterminez l'expression explicite de $\hat{\theta}$ dans ce cas (où ici $\hat{\theta} = \operatorname{Argmin}_{\theta \in \mathbb{R}^p} \tilde{\mathcal{L}}(\theta)$).

Exercice 1.

1. $\hat{f}_{1,ppv}(x_1) = +$ (car le plus proche voisin est affecté à un \oplus).

$\hat{f}_{3,ppv}(x_9) = -$ (car 2 des 3 ppv sont \ominus).

2. Sur 1/2/3/4, on fabrique un prédicteur \hat{f}_1 .
Sur 5/6/7/8, on fabrique un prédicteur \hat{f}_2 .

Comme x est un 3ppv, on a alors : $\hat{f}_1(x) = \ominus \quad \forall x$
 $\hat{f}_2(x) = \oplus \quad \forall x$.

De sorte que \hat{f}_1 se trompe 3 fois sur 4 sur 5/6/7/8.
 \hat{f}_2 _____ sur 1/2/3/4.

$$\Rightarrow \hat{R}_{cv} = \frac{1}{2} \left(\frac{3}{4} + \frac{3}{4} \right) = \frac{3}{4}.$$

3. Points non linéairement séparables, au moins 2 points de chaque côté de l'hyperplan \Rightarrow 4 vecteurs supports au minimum, les 2 mal placés et les 2 que l'on trouve déjà dans une séparation de points linéairement séparables.

4. Non traité.

5. On rappelle que l'hétérogénéité globale est ici :

$$H = 2 \left[N_A (\hat{p}_A (1 - \hat{p}_A)) + \underbrace{(N - N_A)}_{N_B} (\hat{p}_B (1 - \hat{p}_B)) \right]$$

avec A, B constituant les 2 feuilles, N_A le nb d'individus dans la feuille A, \hat{p}_A la proportion de + dans A, \hat{p}_B la proportion de - dans B. Si on fait une coupe horizontale laissant les points 2/4/3/9 dans A, on obtient :

$$\frac{H}{2} = 4 \times \frac{1}{4} \times \frac{3}{4} + 5 \times \frac{3}{5} \times \frac{2}{5} = \frac{3}{4} + \frac{6}{5} = \boxed{\frac{9}{5}}$$

Exercice 2:

(2)

$$\begin{aligned}
 1. \quad R_f &= \mathbb{E}[\phi(Y, f(x))] = \mathbb{E}[\mathbb{E}[\phi(Y, f(x)) | x]] \\
 &= \mathbb{E}[\phi(f(x)) \cdot P(Y=1|x)] + \mathbb{E}[\phi(f(-x)) \cdot P(Y=-1|x)] \\
 &= \mathbb{E}[\phi(f(x)) \eta(x) + \phi(f(-x)) (1-\eta(x))] \\
 &= \mathbb{E}[\psi(x)] \text{ avec } \psi(x) = \phi(f(x)) \eta(x) + \phi(f(-x)) (1-\eta(x)).
 \end{aligned}$$

$$\begin{aligned}
 2. \quad &\text{Comme } f \text{ est à valeurs dans } \{-1, 1\} \\
 &\psi(x) \in \{ \phi(1) \eta(x) + (1-\eta(x)) \phi(-1), \phi(-1) \eta(x) + (1-\eta(x)) \phi(1) \} \\
 \Rightarrow &\psi(x) \geq \min \{ \eta(x)(\phi(1) - \phi(-1)) + \phi(-1), \eta(x)(\phi(-1) - \phi(1)) + \phi(1) \} \\
 \Rightarrow &R_f \geq \mathbb{E}[\min \{ \eta(x)(\phi(1) - \phi(-1)) + \phi(-1), \eta(x)(\phi(-1) - \phi(1)) + \phi(1) \}]
 \end{aligned}$$

Il nous reste à montrer que ce minorant est un minimum.
 Posons $a = -\phi(1) + \phi(-1)$. On constate que $\forall \eta \in [0, 1]$,

$$(*) \quad -a\eta + \phi(-1) \geq a\eta + \phi(1) \Leftrightarrow 2a\eta \leq a \Leftrightarrow \eta \leq 1/2$$

car $a > 0$ par hypothèse. Ainsi, si l'on pose

$$f^*(x) = \begin{cases} 1 & \text{si } \eta(x) \geq 1/2 \\ -1 & \text{si } \eta(x) < 1/2 \end{cases}$$

$$\text{on constate que } \psi(x) = \begin{cases} \phi(1)\eta(x) + \phi(-1)(1-\eta(x)) & \text{si } \eta(x) \geq 1/2 \\ \phi(-1)\eta(x) + \phi(1)(1-\eta(x)) & \text{si } \eta(x) < 1/2 \end{cases}$$

$$(**) \quad \begin{cases} -a\eta(x) + \phi(-1) & \text{si } \eta(x) \geq 1/2 \\ a\eta(x) + \phi(1) & \text{si } \eta(x) < 1/2 \end{cases}$$

$$= \min \{ -a\eta(x) + \phi(-1), a\eta(x) + \phi(1) \}$$

d'après (*). Il vient $R_{f^*} = \mathbb{E}[\min \{ -a\eta(x) + \phi(-1), a\eta(x) + \phi(1) \}]$
 puis le résultat.

3. Si $\eta(x) = x$, alors en posant $a = e - e^{-1}$ (3)

$$\Psi_{f^*}(x) = \begin{cases} ax + e^{-1} & \text{si } x \leq 1/2 \\ -ax + e & \text{si } x \geq 1/2 \end{cases} \quad (\text{d'après (**)})$$

$$\begin{aligned} \text{Ainsi, } R_{f^*} &= \mathbb{E}[\Psi_{f^*}(X)] \\ &= a \mathbb{E}[X \mathbf{1}_{\{X \leq 1/2\}}] + e^{-1} P(X \leq 1/2) \\ &\quad - a \mathbb{E}[X \mathbf{1}_{\{X \geq 1/2\}}] + e P(X \geq 1/2) \end{aligned}$$

$$\mathbb{E}[X \mathbf{1}_{\{X \leq 1/2\}}] = \int_0^{1/2} x dx = \left[\frac{x^2}{2} \right]_0^{1/2} = \frac{1}{8} \quad \text{et} \quad \mathbb{E}[X \mathbf{1}_{\{X \geq 1/2\}}] = \left[\frac{x^2}{2} \right]_{1/2}^1 = \frac{3}{8}$$

de sorte que :

$$R_{f^*} = -\frac{1}{4}a + \frac{e+e^{-1}}{2} = \boxed{-\frac{1}{2} \operatorname{sh}(1) + \operatorname{ch}(1)}.$$

4. (a) Il est naturel de poser :

$$\hat{\Theta}_N = \underset{\Theta \in \Theta}{\operatorname{Argmin}} \quad \frac{1}{N} \sum_{k=1}^N \phi(Y_k \cdot f_{\Theta}(X_k)).$$

car par la LGN, $\frac{1}{N} \sum_{k=1}^N \phi(Y_k \cdot f_{\Theta}(X_k)) \xrightarrow{N \rightarrow \infty} \mathbb{E}[\phi(Y \cdot f_{\Theta}(X))]$
 (si $\mathbb{E}[|\phi(Y \cdot f_{\Theta}(X))|] < \infty$)

(b) Comment calculer $\hat{\Theta}_n$? L'idée proposée ici est une simple descente de gradient de pas γ . En effet, comme

$$\nabla_{\Theta} e^{-y \cdot f_{\Theta}(x)} = -y \nabla_{\Theta} f_{\Theta}(x) e^{-y \cdot f_{\Theta}(x)}, \quad \text{on constate}$$

que (Θ_k) est bien de la forme :

$$\Theta_{k+1} = \Theta_k - \gamma \times \nabla_{\Theta} \left(\frac{1}{N} \sum_{i=1}^N e^{-Y_i f_{\Theta_k}(X_i)} \right)$$

• A propos de la convergence lorsque $f_0(u) = \langle u, \phi \rangle$.

(1)

• la descente de gradient converge vers ϕ^* lorsque

la fonction à minimiser est convexe avec un unique minimum ϕ^* . Ici, le fait que ϕ^* soit l'unique minimum est admis (Pas toujours vrai évidemment). Il nous reste pour vérifier la convergence vers ϕ^* à montrer que $\phi \rightarrow \frac{1}{N} \sum_{k=1}^N e^{-\gamma_k f_0(x_k)}$ est convexe.

Or, $x \rightarrow \langle x, \phi \rangle$ est convexe (car linéaire) et $u \rightarrow e^{\lambda u}$ est convexe $\forall \lambda \in \mathbb{R} \Rightarrow \phi \rightarrow e^{-\gamma f_0(u)}$ est convexe $\forall (u, \gamma)$.

Le résultat suit.

(c) Pour une version mini-batch, il suffit de faire un tirage aléatoire $(I_1^{(k)}, \dots, I_M^{(k)})$ de M observations ($M \ll N$) à chaque itération et d'écrire:

$$\phi_{k+1}^{(NB)} = \phi_k^{(NB)} + \frac{\gamma}{M} \sum_{i=1}^M \nabla_{\phi} f_{\phi_k}^{(NB)}(x_{I_i^{(k)}}) \frac{y_{I_i^{(k)}}}{I_i^{(k)}} e^{-\frac{\gamma}{I_i^{(k)}} f_{\phi_k}(x_{I_i^{(k)}})}$$

(d) Classiquement, si Θ est "grand" (de grande dimension) la convergence de l'algorithme de descente sera plus lente mais surtout, on se confronte à un risque de surapprentissage^(*) lié à une variance trop élevée de $\hat{\phi}_N$. Si Θ est "petit", la variance est plus faible mais le biais est plus important.

(*) : Dans le cadre des réseaux de neurones, on constate néanmoins qu'il n'y a pas autant de surapprentissage que ce que la dimension de Θ devrait générer. Il est "communément" admis que l'algorithme choisit des minimiseurs qui sont de "bonnes solutions" ---

Exercice 3

1) Il s'agit d'un problème de régression pénalisé avec une pénalisation L^1 sur les accroissements. On veut donc faire en sorte que peu d'accroissements soient \neq de 0 (la pénalisation L^1 a vocation à approximer la pénalisation L^0 qui elle même est conçue pour gérer des de la parcimonie).

Application: Données "géographiques" (Météo, Capteurs sur un cerveau...)

2) $\Theta \rightarrow |\Theta_{j+1} - \Theta_j|$ est clairement convexe (par exemple, $\forall \lambda \in [0,1]$,

$$\text{on a bien} \quad |(\lambda \Theta + (1-\lambda) \tilde{\Theta})_{j+1} - (\lambda \Theta + (1-\lambda) \tilde{\Theta})_j| \leq \lambda |\Theta_{j+1} - \Theta_j| + (1-\lambda) |\tilde{\Theta}_{j+1} - \tilde{\Theta}_j|.$$

Ainsi, \mathcal{L} est convexe en tant que somme de fonctions convexes.

Notons $X = (X_1 | \dots | X_p)$, $\forall j \in \{2, \dots, p-1\}$,

$$\partial_{\Theta_j} \mathcal{L} = 2N X_j^T (X\Theta - Y) + \lambda \begin{cases} \text{sgn}(\Theta_{j+1} - \Theta_j) + \text{sgn}(\Theta_j - \Theta_{j-1}) & \text{si } \Theta_j \neq \Theta_{j-1} \text{ et } \Theta_j \neq \Theta_{j+1} \\ \text{sgn}(\Theta_{j+1} - \Theta_j) + [-1, 1] & \text{si } \Theta_j \neq \Theta_{j+1} \text{ et } \Theta_{j-1} = \Theta_j \\ [-1, 1] + \text{sgn}(\Theta_j - \Theta_{j-1}) & \text{si } \Theta_j = \Theta_{j+1} \text{ et } \Theta_{j-1} \neq \Theta_j \\ 2[-1, 1] & \text{si } \Theta_{j-1} = \Theta_j = \Theta_{j+1} \end{cases}$$

3) Si $X^T X = I_p$, alors $X_j^T X \Theta = \Theta_j$.

Si $\Theta_1 < \dots < \Theta_p$, alors, il vient: $\forall j \in [2, p-1]$

$$\partial_{\Theta_j} \mathcal{L}(\hat{\Theta}) = 0 = 2N(\hat{\Theta}_j - X_j^T Y) + 2\lambda \Rightarrow \forall j \in [2, p-1], \quad \hat{\Theta}_j = X_j^T Y - \frac{\lambda}{N}.$$

4) Sous ces conditions on constate que

$$-2N X_j^T Y + 2\lambda [-1, 1] \text{ contient bien } 0 \quad \forall j \in [2, p-1]$$

Pour les cas $j=1$, $j=p$, il y a une petite différence dans le

sous-gradient et il faut donc que:

$$-2N X_j^T Y + \lambda [-1, 1] \text{ contienne } 0 \quad (j'=1 \text{ ou } p). \quad (6)$$

C'est bien le cas d'ailleurs sous les conditions de la question.

Il vient $0 \in \{0, \underset{\substack{\uparrow \\ \text{sous-gradient}}}{\nabla L(0)} = 0\} \Rightarrow 0 \text{ est minimum}$

via un résultat du cours.

$$5. (a) \begin{pmatrix} \theta_2 - \theta_1 \\ \theta_3 - \theta_2 \\ \vdots \\ \theta_p - \theta_{p-1} \\ \theta_p \end{pmatrix} = \underbrace{\begin{pmatrix} -1 & 1 & & & \\ & \ddots & \ddots & & \\ & & -1 & 1 & \\ & & & -1 & 1 \\ & & & & 1 \end{pmatrix}}_D \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_p \end{pmatrix}$$

\Rightarrow résultat.

(b) On a donc à chercher "le" minimum de

$$\tilde{L}(\theta) = \|Y - X\theta\|^2 + \lambda \|D\theta\|_1.$$

$$\beta = D\theta \quad \Rightarrow \quad \|Y - XD^{-1}\beta\|^2 + \lambda \|\beta\|_1$$

(D est clairement inversible).

On a donc un problème de type lasso classique avec

$\tilde{X} = XD^{-1}$ à la place de X. Sous les hypothèses,

$\tilde{X}^T \tilde{X} = I_p$. On sait alors que la solution (voir cours dans une forme légèrement différente en terme de normalisation)

est: $\hat{\beta}^{\text{lasso}} = (\hat{\beta}_j^{\text{lasso}})_{j=1}^p$ avec

$$\hat{\beta}_j^{\text{lasso}} = \text{sgn}(\hat{\beta}_j) \left(|\hat{\beta}_j| - \frac{\lambda}{2N} \right)_+$$

où $\hat{\beta} = \tilde{X}^T Y$. On a alors $\hat{\theta} = D^{-1} \hat{\beta}$.