

Dans la suite, on va noter $\eta(x) = P(Y=1|X=x)$

②

Problème 2: l'ensemble de Bayes (unique associé à la perturbation)

$$R^* = \inf_{g: X \rightarrow \{-1, +1\}} P[g(X) \neq Y].$$

$$g: X \rightarrow \{-1, +1\}$$

Montrons l'équivalence entre les deux définitions de R^* :

$$a) \quad R^* = \mathbb{E}[\min(\eta(X), 1 - \eta(X))]$$

$$b) \quad R^* = \frac{1}{2} - \frac{1}{2} \mathbb{E}[|2\eta(X) - 1|].$$

c) Supposons que pour $y \in \{-1, +1\}$, les variables $(X|Y=y)$ ont f_+ et f_- comme densités respectives pour $y = -1$ et $y = +1$, et

$$P(Y=1) = P(Y=-1) = \frac{1}{2}. \quad \text{Montrons que:}$$

$$R^* = \frac{1}{2} - \frac{1}{4} \int |f_+(x) - f_-(x)| dx.$$

Solution:

$$(*) \quad \text{Rappelons que } g^*(x) = \begin{cases} +1 & \text{si } \eta(x) \geq \frac{1}{2} \\ -1 & \text{sinon} \end{cases}$$

$$\text{i.e. } R^* = R(g^*) = P(g^*(X) \neq Y).$$

Nous avons:

$$P(g^*(X) \neq Y | X=x) = P(g^*(x) = +1, Y = -1 | X=x) + P(g^*(x) = -1, Y = +1 | X=x)$$

$$\begin{aligned}
&= (1 - \eta(x)) \mathbb{1}_{\{g^*(x) = +1\}} + \eta(x) \mathbb{1}_{\{g^*(x) = -1\}} \quad (2) \\
&= (1 - \eta(x)) \mathbb{1}_{\{\eta(x) \geq \frac{1}{2}\}} + \eta(x) \mathbb{1}_{\{\eta(x) < \frac{1}{2}\}} \\
&= (1 - \eta(x)) \mathbb{1}_{\{\eta(x) \geq 1 - \eta(x)\}} + \eta(x) \mathbb{1}_{\{\eta(x) < 1 - \eta(x)\}} \\
&= \min(\eta(x), 1 - \eta(x)).
\end{aligned}$$

Esprance par rapport à la loi de X donne:

$$R^* = \mathbb{E}[\min(\eta(X), 1 - \eta(X))]$$

b) Nous avons:

$$\min(\eta(x), 1 - \eta(x)) = \frac{1}{2} - \left| \eta(x) - \frac{1}{2} \right|$$

$$\left\{ \begin{array}{l} \text{si } \eta(x) \geq 1 - \eta(x) \Leftrightarrow \frac{1}{2} - \left| \eta(x) - \frac{1}{2} \right| = \frac{1}{2} - \left(\eta(x) - \frac{1}{2} \right) = 1 - \eta(x) \\ \text{si } \eta(x) < 1 - \eta(x) \Leftrightarrow \frac{1}{2} - \left| \eta(x) - \frac{1}{2} \right| = \frac{1}{2} - \left(\frac{1}{2} - \eta(x) \right) = \eta(x). \end{array} \right.$$

$$R^* = \mathbb{E}[\min(\eta(X), 1 - \eta(X))] = \mathbb{E}\left[\frac{1}{2} - \left| \eta(X) - \frac{1}{2} \right|\right] = \frac{1}{2} - \frac{1}{2} \mathbb{E}[|2\eta(X) - 1|]$$

(c) Théorème de Bayes:

$$\eta(x) = \mathbb{P}(Y=1 | X=x) = \frac{\mathbb{P}(Y=1) f_+(x)}{\mathbb{P}(Y=1) f_+(x) + \mathbb{P}(Y=-1) f_-(x)} = \frac{f_+(x)}{f_+(x) + f_-(x)}$$

$$\text{donc: } 2\eta(x) - 1 = \frac{2f_+(x)}{f_+(x) + f_-(x)} - 1 = \frac{f_+(x) - f_-(x)}{f_+(x) + f_-(x)}$$

$$R^* = \frac{1}{2} - \frac{1}{2} \int \left| \frac{f_+(x) - f_-(x)}{f_+(x) + f_-(x)} \right| \left(\frac{1}{2} f_+(x) + \frac{1}{2} f_-(x) \right) dx = \frac{1}{2} - \frac{1}{4} \int |f_+(x) - f_-(x)| dx$$

Problème ②:

③

La divergence de Kullback-Leibler entre deux vecteurs aléatoires X_1 et X_2 de densités f_1 et f_2 sur un support commun $S \subseteq \mathbb{R}^d$ est donnée par :

$$D(f_1 \| f_2) = \int_S f_1(x) \log \frac{f_1(x)}{f_2(x)} dx.$$

a) Calculer la divergence de Kullback-Leibler pour $X_i \sim \mathcal{N}(\mu_i, \Sigma_i)$ pour $i=1,2$.

b) Supposons que pour $y \in \{-1, +1\}$, $(X|Y=y) \sim \mathcal{N}(\mu_y, \Sigma)$ et $P(Y=+1) = p$. Montrer que

$$\mathbb{E} \left[\sqrt{\eta(X)(1-\eta(X))} \right] = p(1-p) \exp \left(-\frac{1}{4} D(\mu_{+1} \| \mu_{-1}) \right)$$

où $D(\mu_{+1} \| \mu_{-1})$ est la divergence de Kullback-Leibler entre $\mathcal{N}(\mu_{+1}, \Sigma)$ et $\mathcal{N}(\mu_{-1}, \Sigma)$.

Solution:

(a) Supposons: $X \sim f_1 = \mathcal{N}(\mu_1, \Sigma_1)$, nous avons :

$$\mathbb{E}_{f_1} \left[(X - \mu_1)(X - \mu_1)^T \right] = \Sigma_1 \Rightarrow \Sigma_1^{-1} \mathbb{E}_{f_1} \left[(X - \mu_1)(X - \mu_1)^T \right] = \Sigma_1^{-1} \Sigma_1 = I_d \\ \Rightarrow \mathbb{E}_{f_1} \left[\Sigma_1^{-1} (X - \mu_1)(X - \mu_1)^T \right] = I_d \Rightarrow \text{tr} \left(\mathbb{E}_{f_1} \left[\Sigma_1^{-1} (X - \mu_1)(X - \mu_1)^T \right] \right) = \text{tr}(I_d) = d$$

donc: $\mathbb{E}_{f_1} \left[(X - \mu_1)^T \Sigma_1^{-1} (X - \mu_1) \right] = d.$

De manière similaire ;

$$\mathbb{E}_{f_2} \left[(X - \mu_2)^T \Sigma_2^{-1} (X - \mu_2) \right] = \text{tr}(\Sigma_2^{-1} \Sigma_2)$$

④

Comme: $(X - \mu_2)^T \Sigma_2^{-1} (X - \mu_2) = ((X - \mu_1) + (\mu_1 - \mu_2))^T \Sigma_2^{-1} ((X - \mu_1) + (\mu_1 - \mu_2))$

$$= (X - \mu_1)^T \Sigma_2^{-1} (X - \mu_1) + 2(\mu_1 - \mu_2)^T \Sigma_2^{-1} (X - \mu_1) + (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2)$$

On passe à l'espérance:

$$\mathbb{E}_{f_1} \left[(X - \mu_2)^T \Sigma_2^{-1} (X - \mu_2) \right] = \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2)$$

On note que $D(f_1 \| f_2)$ peut être écrite comme une espérance

$$D(f_1 \| f_2) = \mathbb{E}_{f_1} \left[\log \frac{f_1(X)}{f_2(X)} \right] \quad \text{si } X \sim \mathcal{N}(\mu_1, \Sigma_1)$$

$$f_1 = \mathcal{N}(\mu_1, \Sigma_1)$$

$$\text{et } f_2 = \mathcal{N}(\mu_2, \Sigma_2)$$

$$\mathbb{E}_{f_1} \left[\log \frac{f_1(X)}{f_2(X)} \right] = \mathbb{E}_{f_1} \left[\log \frac{|\Sigma_2|^{1/2} \exp\{-\frac{1}{2}(X - \mu_2)^T \Sigma_2^{-1} (X - \mu_2)\}}{|\Sigma_1|^{1/2} \exp\{-\frac{1}{2}(X - \mu_1)^T \Sigma_1^{-1} (X - \mu_1)\}} \right]$$

$$= \frac{1}{2} \log \frac{|\Sigma_2|}{|\Sigma_1|} + \frac{1}{2} \mathbb{E}_{f_1} \left[(X - \mu_2)^T \Sigma_2^{-1} (X - \mu_2) - (X - \mu_1)^T \Sigma_1^{-1} (X - \mu_1) \right]$$

$$= \frac{1}{2} \left[\log \frac{|\Sigma_2|}{|\Sigma_1|} + \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2) - d \right]$$

(b) On note $f_+ = \mathcal{N}(\mu_+, \Sigma)$ et $f_- = \mathcal{N}(\mu_-, \Sigma)$, alors:

$$f = p f_- + (1-p) f_+; \quad \text{densité marginale en } X.$$

$$\begin{aligned} \mathbb{E} \left(\sqrt{\eta(X)(1-\eta(X))} \right) &= \int \sqrt{\eta(x)(1-\eta(x))} f(x) dx \\ &= \int \sqrt{\frac{p f_-(x)}{f(x)} \frac{(1-p) f_+(x)}{f(x)}} f(x) dx = \sqrt{p(1-p)} \int \sqrt{f_-(x) f_+(x)} dx \\ &= \sqrt{p(1-p)} \int \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{4} h(x)\right\} dx \end{aligned}$$

(5)

$$\text{on a } Q(x) = (x - \mu_+)^T \Sigma^{-1} (x - \mu_+) + (x - \mu_-)^T \Sigma^{-1} (x - \mu_-)$$

$$\text{on note } \bar{\mu} = \frac{(\mu_+ + \mu_-)}{2}$$

$$\begin{aligned} Q(x) &= 2(x - \bar{\mu})^T \Sigma^{-1} (x - \bar{\mu}) + \frac{1}{2}(\mu_+ - \mu_-)^T \Sigma^{-1} (\mu_+ + \mu_-) \\ &= 2(x - \bar{\mu})^T \Sigma^{-1} (x - \bar{\mu}) + D(\mu_+ \parallel \mu_-) \end{aligned}$$

on obtient :

$$\mathbb{E} \left(\sqrt{\gamma(x)(1-\gamma(x))} \right) = \sqrt{p(1-p)} \exp \left\{ -\frac{1}{4} D(\mu_+ \parallel \mu_-) \right\}$$

$$\begin{aligned} &\int \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \bar{\mu})^T \Sigma^{-1} (x - \bar{\mu}) \right\} dx \\ &= \sqrt{p(1-p)} \exp \left\{ -\frac{1}{4} D(\mu_+ \parallel \mu_-) \right\} \end{aligned}$$

Problème ③: Dans cet exercice, nous explorons quelques questions liées aux données de grande dimension.

Pour $i=1, \dots, n$, soit $X^{(i)} \in \mathbb{R}^d$ des vecteurs aléatoires i.i.d de loi uniforme $[0,1]^d$. Soit X une nouvelle observation, on définit la expérience de la distance à la donnée la plus proche

$$f_{\infty}(d, n) = \mathbb{E} \left[\min_{i=1, \dots, n} \|X^{(i)} - X\|_{\infty} \right].$$

L'objectif de cet exercice est de comprendre le comportement de $f_{\infty}(d, n)$ lorsque n et d sont grands.

(a) Montrer que pour tout $t > 0$

$$\mathbb{P} \left[\min_{i=1, \dots, n} \|X^{(i)} - X\|_{\infty} > t \right] \geq 1 - n(2t)^d.$$

(b) Supposons qu'on veut être sûr que X est une distance $\frac{1}{4}$ à au moins une donnée $x^{(1)}, x^{(2)}, \dots, x^{(n)}$ avec une probabilité d'au moins $\frac{1}{2}$. Trouver la borne inférieure de la taille du jeu de données n en fonction de d .

$$\begin{aligned}
 \mathbb{P}\left[\min_{i=1, \dots, n} \|x^{(i)} - X\|_d > t\right] &= \mathbb{P}\left[\bigcap_{i=1}^n \{ \|x^{(i)} - X\|_d > t \}\right] \\
 &= 1 - \mathbb{P}\left[\bigcup_{i=1}^n \{ \|x^{(i)} - X\|_d \leq t \}\right] = 1 - \mathbb{P}\left[\bigcup_{i=1}^n \{ \|x^{(i)} - X\|_d \leq t \}\right] \\
 &= 1 - \mathbb{P}\left(\bigcup_{i=1}^n \{ \|x^{(i)} - X\|_d \leq t \}\right) \\
 &\geq 1 - \sum_{i=1}^n \mathbb{P}(\{ \|x^{(i)} - X\|_d \leq t \}) \\
 &= 1 - n \mathbb{P}(\{ \|x^{(1)} - X\|_d \leq t \}) \\
 &= 1 - n \mathbb{P}\left[\bigcap_{j=1, \dots, d} \{ |x_j^{(1)} - x_j| \leq t \}\right] \\
 &= 1 - n \left(\mathbb{P}[|x_1^{(1)} - x_1| \leq t] \right)^d \\
 &= 1 - n \left(1 - (1-t)^2 \right)^d \\
 &= 1 - n (2t - t^2)^d \\
 &\geq 1 - n (2t)^d \quad \text{si } 0 < t < 1.
 \end{aligned}$$

(b) On s'intéresse à $\mathbb{P} \left[\min_{i=1, \dots, n} \|X^{(i)} - x\|_\infty \leq \frac{1}{4} \right]$. (7)

$$\begin{aligned} \mathbb{P} \left[\min_{i=1, \dots, n} \|X^{(i)} - x\|_\infty \leq \frac{1}{4} \right] &= 1 - \mathbb{P} \left[\min_{i=1, \dots, n} \|X^{(i)} - x\|_\infty > \frac{1}{4} \right] \\ &\leq 1 - \left(1 - n \left(\frac{1}{2} \right)^d \right) \\ &= n \left(\frac{1}{2} \right)^d. \end{aligned}$$

pour garantir que cette probabilité $\geq \frac{1}{2}$, on doit avoir

$$\frac{n}{2^d} \geq \frac{1}{2} \Leftrightarrow n \geq 2^{d-1}.$$

(c) En utilisant (a), montrer que $f(d, n) \geq \frac{d}{2(d+1)} n^{-\frac{1}{d}}$.

Hint: $[Z \geq 0; \mathbb{E}(Z) = \int_0^{+\infty} \mathbb{P}(Z \geq t) dt]$

$$\begin{aligned} f(d, n) &= \int_0^{+\infty} \mathbb{P} \left[\min_{1 \leq i \leq n} \|X^{(i)} - x\|_\infty > t \right] dt \\ &\geq \int_0^{+\infty} \left(1 - n(2t)^d \right)_+ dt = \int_0^{\frac{1}{2} n^{-\frac{1}{d}}} \left(1 - n(2t)^d \right) dt \\ &= \boxed{\frac{d}{2(d+1)} n^{-\frac{1}{d}}} \end{aligned}$$

(d) Calculer cette quantité pour $d \in \{2, 10, 20\}$ et $n \in \{100, 1000, 10000, 100000\}$