

## TD 04 : Régression linéaire

Fondements du Machine Learning, L3 IM2D

Octobre-Novembre 2024 (V2 : 3 novembre 2024)



### Exercice 4.1 : Variation de données

(Adapté d'un exercice d'examen proposé en 2019-2020.)

On considère un jeu de données  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ , où  $\mathbf{x}_i \in \mathbb{R}^n$  et  $y_i \in \mathbb{R}$  pour tout  $i = 1, \dots, m$ . On notera de manière plus compacte

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_m^\top \end{bmatrix} \in \mathbb{R}^{m \times n} \quad \text{et} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^m.$$

On suppose dans cet exercice que  $\text{rang}(\mathbf{X}) = n \leq m$ , et qu'il existe un vecteur  $\boldsymbol{\beta}^* \in \mathbb{R}^n$  tel que:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\epsilon},$$

où  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 I_m)$ , c'est-à-dire que les composantes de  $\boldsymbol{\epsilon}$  sont i.i.d. selon une loi normale de moyenne 0 et de variance  $\sigma^2 > 0$ .

On considère ensuite le problème aux moindres carrés suivant :

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^n} f(\boldsymbol{\beta}) := \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i^\top \boldsymbol{\beta} - y_i)^2 = \frac{1}{m} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2. \quad (1)$$

- Donner une solution du problème (1) sous les hypothèses considérées. Pourquoi cette solution est-elle unique ?
- Soit  $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ . Montrer que  $\mathbf{H}$  est symétrique et que  $\mathbf{H}^2 = \mathbf{H}$ .
- Justifier qu'une valeur propre de  $\mathbf{H}$  est nécessairement égale à 0 ou 1. Sachant que  $\text{rang}(\mathbf{H}) = n$ , en déduire que  $\text{trace}(\mathbf{H}) = n$ .
- Montrer que  $\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{H}\mathbf{y}$ , où  $\hat{\boldsymbol{\beta}}$  est la solution de (1) considérée en question a). Montrer alors que

$$f(\hat{\boldsymbol{\beta}}) = \frac{1}{m} (\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{H} \mathbf{y}).$$

e) Justifier que  $\mathbf{H}\mathbf{X}\boldsymbol{\beta}^* = \mathbf{X}\boldsymbol{\beta}^*$ , où  $\boldsymbol{\beta}^*$  est la vérité terrain du problème. En déduire que :

$$f(\hat{\boldsymbol{\beta}}) = \frac{1}{m}(\boldsymbol{\epsilon}^T \boldsymbol{\epsilon} - \boldsymbol{\epsilon}^T \mathbf{H} \boldsymbol{\epsilon}).$$

f) On veut maintenant obtenir l'espérance de l'erreur représentée par  $f$  lorsque le jeu de données  $\mathcal{D}$  (donc le vecteur  $\boldsymbol{\epsilon}$ ) varie. Montrer que l'on a :

$$\mathbb{E}_{\mathcal{D}} [f(\hat{\boldsymbol{\beta}})] = \sigma^2 \left(1 - \frac{n}{m}\right).$$

*Indication :  $\mathbb{E}_{\mathbf{u}} [\mathbf{u}^T \mathbf{M} \mathbf{u}] = \sigma^2 \text{trace}(\mathbf{M})$  pour toute matrice carrée  $\mathbf{M} \in \mathbb{R}^{n \times n}$  et tout vecteur  $\mathbf{u} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ .*

g) En déduire que

$$\mathbb{E}_{\mathcal{D}} [f(\hat{\boldsymbol{\beta}})] \longrightarrow \mathbb{E}_{(x,y)} [R(\boldsymbol{\beta}^*)] \text{ quand } m \rightarrow \infty,$$

où  $R(\boldsymbol{\beta}) = (\mathbf{x}^T \boldsymbol{\beta} - y)^2$  et  $(x, y)$  suit la même distribution que les éléments de  $\mathcal{D}$ , c'est-à-dire que l'on a  $y = \mathbf{x}^T \boldsymbol{\beta}^* + \epsilon$  avec  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ .

### Exercice 4.2 : Ajout de données

Soient  $\mathbf{X} \in \mathbb{R}^{m \times n}$  et  $\mathbf{y} \in \mathbb{R}^m$  avec  $\text{rang}(\mathbf{X}) = n \leq m$  et  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\epsilon}$ , où  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_m)$ . On rappelle que dans ce contexte, l'estimateur du maximum de vraisemblance (ou l'estimateur des moindres carrés) vaut

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \boldsymbol{\beta}^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}.$$

a) Retrouver les formules  $\mathbb{E} [\hat{\boldsymbol{\beta}}] = \boldsymbol{\beta}^*$  ainsi que

$$\text{Cov} [\hat{\boldsymbol{\beta}}] = \mathbb{E} [(\hat{\boldsymbol{\beta}} - \mathbb{E} [\hat{\boldsymbol{\beta}}])(\hat{\boldsymbol{\beta}} - \mathbb{E} [\hat{\boldsymbol{\beta}}])^T] = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}.$$

b) Pour tout vecteur aléatoire  $\mathbf{z} \in \mathbb{R}^n$ , montrer que

$$\mathbb{E} [\|\mathbf{z} - \mathbb{E} [\mathbf{z}]\|_2^2] = \text{trace Cov} [\mathbf{z}].$$

En déduire que  $\mathbb{E} [\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_2^2] = \sigma^2 \text{trace} ((\mathbf{X}^T \mathbf{X})^{-1})$ .

c) On considère l'ajout d'une nouvelle donnée  $(\mathbf{v}, w = \mathbf{v}^T \boldsymbol{\beta}^* + \epsilon_{m+1}) \in \mathbb{R}^n \times \mathbb{R}$  avec  $\epsilon_{m+1} \sim \mathcal{N}(0, \sigma^2)$ , et on s'intéresse à l'impact de cet ajout sur l'estimateur du maximum de vraisemblance. Etablir la formule de cet estimateur pour le nouveau jeu de données (on le notera  $\tilde{\boldsymbol{\beta}}$ ) en fonction de  $\mathbf{X}, \mathbf{v}, \mathbf{y}, w$  et calculer sa matrice de covariance.

d) Pour toute matrice  $\mathbf{A} \in \mathbb{R}^{n \times n}$  définie positive et tout vecteur  $\mathbf{u} \in \mathbb{R}^n$ , on a :

$$(\mathbf{A}^{-1} + \mathbf{u}\mathbf{u}^T)^{-1} = \mathbf{A} - \frac{1}{1 + \mathbf{u}^T \mathbf{A} \mathbf{u}} \mathbf{A} \mathbf{u} \mathbf{u}^T \mathbf{A}.$$

En utilisant cette relation, montrer que  $\text{trace}((\mathbf{A}^{-1} + \mathbf{u}\mathbf{u}^T)^{-1}) = \text{trace}(\mathbf{A}) - \frac{\mathbf{u}^T \mathbf{A} \mathbf{u}}{1 + \mathbf{u}^T \mathbf{A} \mathbf{u}}$ .

e) Déduire de la question précédente un lien entre  $\mathbb{E} [\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_2^2]$  et  $\mathbb{E} [\|\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_2^2]$ . Sous quelle condition est-il avantageux d'introduire une donnée supplémentaire ?

### Exercice 4.3 : Maximum de vraisemblance

(Adapté d'un exercice d'examen proposé en 2020-2021.)

On considère un jeu de données sous la forme de deux matrices  $\mathbf{X} \in \mathbb{R}^{m_1 \times n}$ ,  $\mathbf{W} \in \mathbb{R}^{m_2 \times n}$ , et de deux vecteurs  $\mathbf{y} \in \mathbb{R}^{m_1}$ ,  $\mathbf{z} \in \mathbb{R}^{m_2}$ . On fait l'hypothèse d'un modèle linéaire sous-jacent affecté par un bruit gaussien, c'est-à-dire que l'on suppose qu'il existe deux vecteurs  $\beta_1^* \in \mathbb{R}^n$  et  $\beta_2^* \in \mathbb{R}^n$  tels que

$$\mathbf{y} = \mathbf{X}\beta_1^* + \epsilon_1 \quad \text{et} \quad \mathbf{z} = \mathbf{W}\beta_2^* + \epsilon_2,$$

où  $\epsilon_1$  est un vecteur aléatoire de  $\mathbb{R}^{m_1}$  distribué selon une loi gaussienne  $\mathcal{N}(\mathbf{0}_{\mathbb{R}^{m_1}}, \sigma^2 \mathbf{I}_{m_1})$  et  $\epsilon_2$  est un vecteur de  $\mathbb{R}^{m_2}$  distribué selon une loi gaussienne  $\mathcal{N}(\mathbf{0}_{\mathbb{R}^{m_2}}, \zeta^2 \mathbf{I}_{m_2})$ , avec  $\sigma > 0$  et  $\zeta > 0$ . On supposera également que les matrices  $\mathbf{X}$  et  $\mathbf{W}$  sont de même rang  $n$ , avec  $n \leq \min\{m_1, m_2\}$ .

- a) Quelles sont alors les distributions des vecteurs  $\mathbf{y}$  et  $\mathbf{z}$  ?
- b) Rappeler la définition d'une fonction de vraisemblance pour  $m$  observations réelles  $a_1, \dots, a_m$  dépendant d'un vecteur de paramètres  $\beta$ , notée  $L(a_1, \dots, a_m; \beta)$ . En suivant cette définition, établir des formules pour les vraisemblances

$$L(y_1, \dots, y_{m_1}; \beta_1) \quad \text{et} \quad L(z_1, \dots, z_{m_2}; \beta_2)$$

- c) On s'intéresse à l'estimateur du maximum de vraisemblance pour  $\beta_1^*$ .
- Donner le problème de maximisation de la vraisemblance dont cet estimateur est solution. Pourquoi est-il équivalent de maximiser la vraisemblance et son logarithme ?
  - Donner le lien entre l'estimateur du maximum de vraisemblance et l'ensemble des solutions du problème aux moindres carrés

$$\min_{\beta_1 \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{X}\beta_1 - \mathbf{y}\|^2.$$

- En déduire une formule pour l'estimateur du maximum de vraisemblance de  $\beta_1^*$ , que l'on notera  $\hat{\beta}_1$ .
- d) En utilisant le même raisonnement qu'en question c), donner la formule de l'estimateur du maximum de vraisemblance pour  $\beta_2^*$ , noté  $\hat{\beta}_2$ .

- e) Justifier que le vecteur  $\begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix}$  est un estimateur non biaisé du vecteur  $\begin{bmatrix} \beta_1^* \\ \beta_2^* \end{bmatrix}$ .

## Exercice 4.4 : Maximum a posteriori

(Adapté d'un exercice d'examen proposé en 2019-2020.)

Soit un jeu de données  $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$ . Dans le cadre de la régression linéaire, on recherche un vecteur  $\beta$  tel que

$$y_i \approx \mathbf{x}_i^T \beta + \epsilon_i,$$

où les  $\epsilon_i$  sont i.i.d. selon une loi normale de moyenne 0 et de variance 1.

On s'intéresse ici à l'estimateur du maximum a posteriori obtenu avec un a priori uniforme sur un intervalle réel  $[a, b]$  (avec  $a < b$ ). On rappelle que la densité de probabilité d'une variable  $z$  suivant une loi uniforme dans un intervalle  $[a, b]$  est la fonction  $p$  donnée par :

$$p(z) = \begin{cases} \frac{1}{b-a} & \text{si } z \in [a, b], \\ 0 & \text{sinon.} \end{cases} \quad (2)$$

- Rappeler la définition de la vraisemblance  $\mathcal{L}(y_1, \dots, y_m; \beta)$ , et donner sa formule pour le problème considéré.
- On suppose a priori que les entrées de  $\beta$  sont i.i.d. suivant une loi uniforme sur  $[a, b]$ . Donner alors la formule de la loi a priori de  $\beta$ , notée  $\mathcal{L}(\beta)$ .
- Formuler le problème d'estimation du maximum a posteriori en utilisant les logarithmes des lois de probabilité.
- Montrer que si l'on ne tient pas compte des termes constants, le problème se ramène à celui ci-dessous :

$$\max_{\beta \in \mathbb{R}^n} \left\{ -\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^T \beta - y_i)^2 + \sum_{j=1}^n \ln(\chi_{b-a}([\beta]_j)) \right\}, \quad (3)$$

où  $\chi_{b-a}(z) = 1$  si  $a \leq z \leq b$  et  $\chi_{b-a}(z) = 0$  sinon.

- On peut montrer que l'ensemble des solutions du problème (3) correspond à celui du problème aux moindres carrés avec contraintes suivant :

$$\min_{\beta \in \mathbb{R}^n} \frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^T \beta - y_i)^2 \quad \text{sous les contraintes } a \leq \beta_j \leq b \quad \forall j = 1, \dots, n. \quad (4)$$

En quoi ce nouveau problème modélise-t-il l'a priori que l'on a mis sur le modèle linéaire ?

- On note  $\hat{\beta}$  l'estimateur du maximum de vraisemblance. Si cet estimateur vérifie  $a \leq [\hat{\beta}]_j \leq b$  pour tout  $j = 1, \dots, n$ , est-il solution du problème (4) ?

## Solutions

### Solution de l'exercice 4.1: Variation de données

a) Une solution du problème est donnée par  $X^\dagger y$ , qui est ici égale à  $(X^T X)^{-1} X^T y$  car on a  $\text{rang}(X) = n$ . On sait dans ce cas que la solution est unique.

b) Un calcul explicite donne

$$H^T = (X(X^T X)^{-1} X^T)^T = (X^T)^T ((X^T X)^{-1})^T X^T = X ((X^T X)^T)^{-1} X^T = X(X^T X)^{-1} X^T,$$

ce qui prouve que  $H$  est symétrique. On vérifie également que

$$H^2 = X(X^T X)^{-1} X^T X(X^T X)^{-1} X^T = X(X^T X)^{-1} X^T = H.$$

c) Soit  $\lambda \in \mathbb{R}$  une valeur propre de  $H$ , et  $v$  un vecteur propre associé. Alors  $Hx = \lambda x$  et  $H^2 x = \lambda^2 x$ . Comme  $H^2 = H$  d'après la question précédente, on en déduit que  $\lambda^2 = \lambda$ , d'où  $\lambda \in \{0, 1\}$ . Comme le rang de  $H$  est égal à  $n$ , cette matrice possède exactement  $n$  valeurs propres non nulles. On en déduit que la trace de  $H$ , qui est égale à la somme des valeurs propres, vaut  $n \times 1 + (m - n) \times 0 = n$ .

d) En utilisant l'expression  $\hat{\beta} = (X^T X)^{-1} X^T y$ , il vient

$$X\hat{\beta} = X(X^T X)^{-1} X^T y = Hy,$$

En utilisant cela dans l'expression de  $f(\hat{\beta})$ , il vient alors

$$\begin{aligned} f(\hat{\beta}) &= \frac{1}{m} \|X\hat{\beta} - y\|^2 \\ &= \frac{1}{m} \|Hy - y\|^2 \\ &= \frac{1}{m} (y^T H^2 y - 2y^T H y + y^T y), \end{aligned}$$

où l'on a utilisé la symétrie de  $H$  pour obtenir le premier et le second termes. En utilisant ensuite que  $H^2 = H$ , on obtient que

$$\begin{aligned} f(\hat{\beta}) &= \frac{1}{m} (y^T H y - 2y^T H y + y^T y) \\ &= \frac{1}{m} (y^T y - y^T H y). \end{aligned}$$

e) On tire de l'expression de  $H$  que

$$\begin{aligned} H X \beta^* &= X(X^T X)^{-1} X^T X \beta^* \\ &= X \beta^*. \end{aligned}$$

Par conséquent, on a  $X\beta^* = H X \beta^*$ . On a ainsi

$$\begin{aligned} y^T y &= (X\beta^* + \epsilon)^T (X\beta^* + \epsilon) \\ &= (X\beta^*)^T X\beta^* + 2\epsilon^T X\beta^* + \epsilon^T \epsilon. \end{aligned}$$

et

$$\begin{aligned} \mathbf{y}^T \mathbf{H} \mathbf{y} &= (\mathbf{X} \boldsymbol{\beta}^* + \boldsymbol{\epsilon})^T \mathbf{H} (\mathbf{X} \boldsymbol{\beta}^* + \boldsymbol{\epsilon}) \\ &= (\mathbf{X} \boldsymbol{\beta}^*)^T \mathbf{H} \mathbf{X} \boldsymbol{\beta}^* + 2 \boldsymbol{\epsilon}^T \mathbf{H} \mathbf{X} \boldsymbol{\beta}^* + \boldsymbol{\epsilon}^T \mathbf{H} \boldsymbol{\epsilon} \\ &= (\mathbf{X} \boldsymbol{\beta}^*)^T \mathbf{X} \boldsymbol{\beta}^* + 2 \boldsymbol{\epsilon}^T \mathbf{X} \boldsymbol{\beta}^* + \boldsymbol{\epsilon}^T \mathbf{H} \boldsymbol{\epsilon} \end{aligned}$$

où l'on a utilisé la symétrie de  $\mathbf{H}$  puis  $\mathbf{H} \mathbf{X} \boldsymbol{\beta}^* = \mathbf{X} \boldsymbol{\beta}^*$ .

En appliquant cela au résultat de la question précédente, on obtient finalement

$$\begin{aligned} f(\hat{\boldsymbol{\beta}}) &= \frac{1}{m} (\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{H} \mathbf{y}) \\ &= \frac{1}{m} ((\mathbf{X} \boldsymbol{\beta}^*)^T \mathbf{X} \boldsymbol{\beta}^* + 2 \boldsymbol{\epsilon}^T \mathbf{X} \boldsymbol{\beta}^* + \boldsymbol{\epsilon}^T \boldsymbol{\epsilon} \\ &\quad - (\mathbf{X} \boldsymbol{\beta}^*)^T \mathbf{X} \boldsymbol{\beta}^* - 2 \boldsymbol{\epsilon}^T \mathbf{X} \boldsymbol{\beta}^* - \boldsymbol{\epsilon}^T \mathbf{H} \boldsymbol{\epsilon}) \\ &= \frac{1}{m} (\boldsymbol{\epsilon}^T \boldsymbol{\epsilon} - \boldsymbol{\epsilon}^T \mathbf{H} \boldsymbol{\epsilon}). \end{aligned}$$

f) On note que la source d'aléatoire provient du vecteur  $\boldsymbol{\epsilon}$ , et donc que  $\mathbb{E}_{\mathcal{D}} [f(\hat{\boldsymbol{\beta}})] = \mathbb{E}_{\boldsymbol{\epsilon}} [f(\hat{\boldsymbol{\beta}})]$ .

En utilisant l'indication et la linéarité de l'espérance, on a

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\epsilon}} [f(\hat{\boldsymbol{\beta}})] &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\epsilon}} [\boldsymbol{\epsilon}^T \boldsymbol{\epsilon}] - \frac{1}{m} \mathbb{E}_{\boldsymbol{\epsilon}} [\boldsymbol{\epsilon}^T \mathbf{H} \boldsymbol{\epsilon}] \\ &= \frac{1}{m} \text{trace}(\sigma^2 \mathbf{I}_m) - \frac{1}{m} \text{trace}(\sigma^2 \mathbf{H}) \\ &= \frac{1}{m} m \sigma^2 - \frac{1}{m} n \sigma^2 = \sigma^2 \left(1 - \frac{n}{m}\right). \end{aligned}$$

g) On tire de la question précédente que  $\mathbb{E}_{\mathcal{D}} [f(\hat{\boldsymbol{\beta}})] \rightarrow \sigma^2$  lorsque  $m \rightarrow \infty$  (mais que  $n$  reste fixe !). Par ailleurs, on a

$$\begin{aligned} \mathbb{E}_{(x,y)} [(x^T \boldsymbol{\beta}^* - y)^2] &= \mathbb{E}_{(x,y)} [(x^T \boldsymbol{\beta}^* - x^T \boldsymbol{\beta}^* - \epsilon)^2] \\ &= \mathbb{E}_{\boldsymbol{\epsilon}} [(\epsilon^2)] = \sigma^2, \end{aligned}$$

où la dernière égalité provient de la définition même de la variance.

## Solution de l'exercice 4.2: Ajout de données

a) On note en préambule que pour tout vecteur aléatoire  $\mathbf{z} \in \mathbb{R}^n$  et toute matrice  $\mathbf{A} \in \mathbb{R}^{n \times n}$  qui ne dépend pas de  $\mathbf{z}$ , on a :

$$\mathbb{E}[\mathbf{A} \mathbf{z}] = \mathbf{A} \mathbb{E}[\mathbf{z}] \quad \text{et} \quad \mathbb{E}[\mathbf{z}^T \mathbf{A}] = \mathbb{E}[\mathbf{z}^T] \mathbf{A}.$$

(Dans le cas d'un vecteur aléatoire, l'espérance est linéaire à gauche et à droite.)

On étudie d'abord l'espérance de  $\hat{\beta}$ . On a :

$$\begin{aligned}
 \mathbb{E} [\hat{\beta}] &= \mathbb{E} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}] \\
 &= \mathbb{E} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \beta^* + \epsilon)] \\
 &= \mathbb{E} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon] \\
 &= \mathbb{E} [\beta^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon] \\
 &= \mathbb{E} [\beta^*] + \mathbb{E} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon] \\
 &= \mathbb{E} [\beta^*] + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E} [\epsilon] \\
 &= \beta^* + 0 = \beta^*.
 \end{aligned}$$

Pour la matrice de covariance, on introduit à nouveau la définition de  $\mathbf{y}$  impliquant  $\beta^*$ . Cela donne :

$$\begin{aligned}
 \mathbb{E} [(\hat{\beta} - \mathbb{E} [\hat{\beta}])(\hat{\beta} - \mathbb{E} [\hat{\beta}])^T] &= \mathbb{E} [(\hat{\beta} - \beta^*)(\hat{\beta} - \beta^*)^T] \\
 &= \mathbb{E} [(\beta^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon - \beta^*)(\beta^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon - \beta^*)^T] \\
 &= \mathbb{E} [((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)^T] \\
 &= \mathbb{E} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}] \\
 &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E} [\epsilon \epsilon^T] \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1},
 \end{aligned}$$

où l'on a utilisé le fait que  $\mathbf{X}^T \mathbf{X} = (\mathbf{X}^T \mathbf{X})^T$ , et donc l'inverse de la transposée de  $(\mathbf{X}^T \mathbf{X})^T$  est égale à l'inverse de  $\mathbf{X}^T \mathbf{X}$ .

Par définition,  $\mathbb{E} [\epsilon] = 0$  et  $\text{Cov} [\epsilon] = I_m$ . Comme  $\text{Cov} [\epsilon] = \mathbb{E} [\epsilon \epsilon^T]$ , on obtient :

$$\begin{aligned}
 \mathbb{E} [(\hat{\beta} - \mathbb{E} [\hat{\beta}])(\hat{\beta} - \mathbb{E} [\hat{\beta}])^T] &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E} [\epsilon \epsilon^T] \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\
 &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = (\mathbf{X}^T \mathbf{X})^{-1},
 \end{aligned}$$

d'où le résultat.

b) On utilise le fait que pour toutes matrices  $\mathbf{A} \in \mathbb{R}^{n \times 1}$  et  $\mathbf{B} \in \mathbb{R}^{1 \times n}$ ,  $\text{trace} \mathbf{BA} = \mathbf{BA} = \text{trace} \mathbf{AB}$ . En effet,

$$\begin{aligned}
 \mathbb{E} [\|z - \mathbb{E} [z]\|_2^2] &= \mathbb{E} [(z - \mathbb{E} [z])^T (z - \mathbb{E} [z])] \\
 &= \mathbb{E} [\text{trace} (z - \mathbb{E} [z])^T (z - \mathbb{E} [z])] \\
 &= \mathbb{E} [\text{trace} (z - \mathbb{E} [z])(z - \mathbb{E} [z])^T] \\
 &= \text{trace} \mathbb{E} [(z - \mathbb{E} [z])(z - \mathbb{E} [z])^T] = \text{trace} \text{Cov} [z],
 \end{aligned}$$

où l'on peut passer de l'avant-dernière à la dernière ligne parce que la trace est une application linéaire (on exploite encore la linéarité de l'espérance).

En appliquant la relation ci-dessus avec  $z = \hat{\beta}$ , on obtient :

$$\mathbb{E} [\|\hat{\beta} - \beta^*\|_2^2] = \mathbb{E} [\|\hat{\beta} - \mathbb{E} [\hat{\beta}]\|_2^2] = \text{trace} (\mathbf{X}^T \mathbf{X})^{-1}.$$

c) La matrice  $\tilde{\mathbf{X}} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_m^T \\ \mathbf{v}^T \end{bmatrix}$  est de rang  $n$ , comme  $\mathbf{X}$ . L'estimateur du maximum de vraisemblance

$$\text{est donc donné par } \tilde{\boldsymbol{\beta}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{y}}, \text{ avec } \tilde{\mathbf{y}} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \\ w \end{bmatrix}.$$

d) On utilise la linéarité de la trace ainsi que la propriété de ce même opérateur déjà évoquée en question b). On a ainsi :

$$\begin{aligned} \text{trace}(\mathbf{A} + \mathbf{u}\mathbf{u}^T)^{-1} &= \text{trace} \mathbf{A} - \frac{1}{1 + \mathbf{u}^T \mathbf{A} \mathbf{u}} \mathbf{A} \mathbf{u} \mathbf{u}^T \mathbf{A} \\ &= \text{trace} \mathbf{A} - \frac{1}{1 + \mathbf{u}^T \mathbf{A} \mathbf{u}} \text{trace} \mathbf{A} \mathbf{u} \mathbf{u}^T \mathbf{A} \\ &= \text{trace} \mathbf{A} - \frac{1}{1 + \mathbf{u}^T \mathbf{A} \mathbf{u}} \text{trace} \mathbf{u}^T \mathbf{A} \mathbf{A} \mathbf{u} \\ &= \text{trace} \mathbf{A} - \frac{1}{1 + \mathbf{u}^T \mathbf{A} \mathbf{u}} \mathbf{u}^T \mathbf{A} \mathbf{A} \mathbf{u}. \end{aligned}$$

e) D'après les questions c) et d), on a :

$$\begin{aligned} \mathbb{E} \left[ \|\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_2^2 \right] &= \text{trace Cov} \left[ \tilde{\boldsymbol{\beta}} \right] \\ &= \text{trace} (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \\ &= \text{trace} (\mathbf{X}^T \mathbf{X} + \mathbf{v}\mathbf{v}^T)^{-1} \\ &= \text{trace} (\mathbf{X}^T \mathbf{X})^{-1} - \frac{\mathbf{v}^T (\mathbf{X}^T \mathbf{X})^{-2} \mathbf{v}}{1 + \mathbf{v}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{v}} \\ &= \mathbb{E} \left[ \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_2^2 \right] - \frac{\mathbf{v}^T (\mathbf{X}^T \mathbf{X})^{-2} \mathbf{v}}{1 + \mathbf{v}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{v}} \end{aligned}$$

Par conséquent,  $\tilde{\boldsymbol{\beta}}$  est un meilleur estimateur (en moyenne quadratique) que  $\hat{\boldsymbol{\beta}}$  lorsque  $\mathbf{v}$  est choisi tel que  $\mathbf{v}^T (\mathbf{X}^T \mathbf{X})^{-2} \mathbf{v}$ , c'est-à-dire  $\|\mathbf{v}\| \neq 0$ .

### Solution de l'exercice 4.3: Maximum de vraisemblance

Le but de l'exercice est de réviser les concepts fondamentaux du maximum de vraisemblance à travers l'utilisation d'un jeu de données double. On pourra rappeler les formules concernant la loi d'une variable gaussienne.

a)  $\mathbf{y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}_1^*, \sigma^2 \mathbf{I}_{m_1})$  et  $\mathbf{z} \sim \mathcal{N}(\mathbf{W}\boldsymbol{\beta}_2^*, \zeta^2 \mathbf{I}_{m_2})$ .

b) On a

$$\mathcal{L}(y_1, \dots, y_{m_1}; \boldsymbol{\beta}_1) = \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^{m_1} \exp \left[ -\frac{1}{2\sigma^2} \|\mathbf{X}\boldsymbol{\beta}_1 - \mathbf{y}\|^2 \right]$$

et

$$\mathcal{L}(z_1, \dots, z_{m_2}; \beta_2) = \left( \frac{1}{\sqrt{2\pi\zeta}} \right)^{m_1} \exp \left[ -\frac{1}{2\zeta^2} \|\mathbf{W}\beta_2 - \mathbf{z}\|^2 \right]$$

c)

i) Le problème de maximisation de la vraisemblance peut être écrit comme

$$\max_{\beta_1 \in \mathbb{R}^n} L(y_1, \dots, y_{m_1}; \beta_1)$$

soit, de manière équivalente par croissance du logarithme, comme

$$\max_{\beta_1 \in \mathbb{R}^n} \ln(L(y_1, \dots, y_{m_1}; \beta_1)).$$

ii) L'estimateur du maximum de vraisemblance est solution du problème aux moindres carrés, car les deux problèmes ont le même ensemble de solutions.

iii) Dans le cadre de l'exercice, comme  $\text{rang}(\mathbf{X}) = n$ , on sait que la solution du problème est unique et égale à  $\hat{\beta}_1 = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ .

d) Par analogie avec la question précédente, on a  $\hat{\beta}_2 = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \mathbf{z}$ .

e) Il suffit de remarquer que  $\mathbb{E}_{\epsilon_1, \epsilon_2} [\hat{\beta}_1] = \beta_1^*$  et  $\mathbb{E}_{\epsilon_1, \epsilon_2} [\hat{\beta}_2] = \beta_2^*$ , puis d'utiliser la formule sur l'espérance d'un vecteur.

### Solution de l'exercice 4.4: Maximum a posteriori

a) La vraisemblance est donnée par la loi jointe des  $y_i$  étant donné le vecteur  $\beta$ . Par hypothèse, on considère que  $y_i - \mathbf{x}_i^T \beta$  suit une loi normale de moyenne 0 et de variance 1, et que les  $y_i$  sont indépendants. Cela conduit à la vraisemblance suivante :

$$\mathcal{L}(y_1, \dots, y_{m_1}; \beta) = \left( \frac{1}{\sqrt{2\pi}} \right)^m \exp \left[ -\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^T \beta - y_i)^2 \right].$$

b) Sous l'hypothèse de la question, la loi a priori de  $\beta$  est donnée par

$$\mathcal{L}(\beta) = \begin{cases} \frac{1}{(b-a)^n} & \text{si } \beta_j \in [a, b] \forall j = 1, \dots, n \\ 0 & \text{sinon.} \end{cases}$$

Par anticipation de la question d), on écrira

$$\mathcal{L}(\beta) = \frac{1}{(b-a)^n} \prod_{j=1}^n \chi_{b-a}([\beta]_j)$$

c) Le problème d'estimation du maximum a posteriori s'écrit

$$\max_{\beta \in \mathbb{R}^n} \ln(\mathcal{L}(y_1, \dots, y_m; \beta) \mathcal{L}(\beta)). \quad (5)$$

d) En utilisant les formules obtenues en questions a) et b), on obtient

$$\begin{aligned} \ln(\mathcal{L}(y_1, \dots, y_m; \boldsymbol{\beta})\mathcal{L}(\boldsymbol{\beta})) &= \ln(\mathcal{L}(y_1, \dots, y_m; \boldsymbol{\beta})) + \ln(\mathcal{L}(\boldsymbol{\beta})) \\ &= m \ln\left(\frac{1}{\sqrt{2\pi}}\right) - \frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^T \boldsymbol{\beta} - y_i)^2 + n \ln\left(\frac{1}{b-a}\right) + \sum_{j=1}^n \ln(\chi_{b-a}([\boldsymbol{\beta}]_j)). \end{aligned}$$

Les termes ne dépendant pas de  $\boldsymbol{\beta}$  n'ont pas d'influence sur la solution du problème d'optimisation (5), on peut les omettre dans la formulation du problème, qui devient donc

$$\max_{\boldsymbol{\beta} \in \mathbb{R}^n} \left\{ -\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^T \boldsymbol{\beta} - y_i)^2 + \sum_{j=1}^n \ln(\chi_{b-a}([\boldsymbol{\beta}]_j)) \right\},$$

soit le résultat attendu.

e) La formulation (4) contraint les valeurs des coefficients de  $\boldsymbol{\beta}$  dans l'intervalle  $[a, b]$ . En supposant un a priori uniformément distribué dans  $[a, b]^n$ , c'est bien ce que l'on cherche à faire.

f) On sait que  $\hat{\boldsymbol{\beta}}$  est la solution du problème

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^n} \frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^T \boldsymbol{\beta} - y_i)^2$$

sans contraintes. Si cet estimateur vérifie de plus  $a \leq [\hat{\boldsymbol{\beta}}]_j \leq b$ , alors il s'agit aussi de la solution du problème (4). En ce sens, le maximum de vraisemblance est compatible avec l'a priori.