

## 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  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_m)$ , 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) ?