

TD 03 : Premiers pas avec le modèle linéaire

Fondements du Machine Learning, L3 IM2D

Octobre 2024



Exercice 3.1 : Modélisation par moindres carrés

On dispose de m notes y_1, \dots, y_m données par des individus ayant chacun passé une nuit dans le même hôtel. On cherche à déterminer une note globale de l'hôtel qui soit la plus cohérente possible avec l'ensemble des m notes.

- Modéliser le problème sous la forme d'un système linéaire et d'un problème aux moindres carrés.
- Pour les deux formulations obtenues, décrire si le problème possède ou non une solution, et donner l'ensemble des solutions s'il existe.
- On suppose maintenant que chaque individu i a passé x_i nuits dans l'hôtel. Comment peut-on tenir compte de cela dans la modélisation ? Quel serait alors l'impact sur la résolution des problèmes associés ?

Exercice 3.2 : Matrices de rang 1

Soit une matrice $\mathbf{X} \in \mathbb{R}^{m \times 1}$ de rang 1.

- Donner une décomposition en valeurs singulières de \mathbf{X} .
- Soit $\mathbf{y} \in \mathbb{R}^m$. Résoudre le problème aux moindres carrés

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

- On suppose que \mathbf{X} est une matrice orthogonale. Que représente alors $\mathbf{X}\beta^*$, où β^* est une solution du problème aux moindres carrés ?

Exercice 3.3 : Moindres carrés régularisés

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

Étant donnés $\mathbf{X} \in \mathbb{R}^{m \times n}$ et $\mathbf{y} \in \mathbb{R}^m$, on considère donc le problème aux moindres carrés suivant :

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\beta}\|_2^2, \quad (1)$$

avec $\lambda \geq 0$ (on englobe ainsi le cas non régularisé vu en cours, qui correspond à $\lambda = 0$).

a) On suppose dans cette question que \mathbf{X} est de rang n avec $n < m$.

- i) Quelle propriété sur $\mathbf{X}^T \mathbf{X}$ cela implique-t-il ?
En déduire la solution de (1) lorsque $\lambda = 0$.
- ii) Lorsque $\lambda > 0$, $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ est définie positive. On peut alors montrer que les solutions du problème (1) sont exactement celles du problème

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^n} \frac{1}{2} \|(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{1/2} \boldsymbol{\beta} - \mathbf{z}\|_2^2,$$

où \mathbf{I} est la matrice identité de $\mathbb{R}^{n \times n}$, et \mathbf{z} vérifie $(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{1/2} \mathbf{z} = \mathbf{X}^T \mathbf{y}$.¹
En déduire l'ensemble des solutions du problème (1) dans ce cas.

b) On suppose maintenant que \mathbf{X} est de rang $m < n$.

- i) Une solution au problème (1) dans le cas $\lambda = 0$ est donné par $\mathbf{X}^\dagger \mathbf{y}$, où $\mathbf{X}^\dagger = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1}$ est la pseudo-inverse de \mathbf{X} . Quelle propriété possède cette solution?
- ii) Lorsque $\lambda > 0$, montrer en utilisant la décomposition en valeurs singulières de \mathbf{X} que la matrice $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ est de rang plein.
- iii) Comment caractériser la ou les solution(s) du problème (1) en utilisant le résultat des questions a)ii) et b)ii) ?

¹On rappelle la notion de racine carrée de matrice déjà utilisée dans le TD 1. Si $\mathbf{A} \succ \mathbf{0}$, alors $\mathbf{A}^{1/2}$ est l'unique matrice symétrique définie positive telle que $(\mathbf{A}^{1/2})^2 = \mathbf{A}$.

Solutions

Solution de l'exercice 3.1: Modélisation par moindres carrés

On peut commencer l'exercice en demandant aux étudiants la note intuitive qu'ils mettraient à l'hôtel. On s'attend à ce que la moyenne soit évoquée, et l'exercice consiste à justifier cela mathématiquement.

a) On cherche une note globale de l'hôtel à partir des y_i . On peut modéliser cela de manière simpliste comme la recherche d'un réel β qui vérifie le système linéaire suivant :

$$\begin{cases} \beta = y_1 \\ \vdots \\ \beta = y_m. \end{cases}$$

On peut aussi chercher la valeur de β qui minimise l'erreur d'approximation des y_i par β . Cela revient à considérer le problème aux moindres carrés

$$\min_{\beta \in \mathbb{R}} \frac{1}{2} \sum_{i=1}^m (\beta - y_i)^2 = \frac{1}{2} \|\mathbf{X}\beta - \mathbf{y}\|^2, \quad \mathbf{X} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \in \mathbb{R}^{m \times 1} \equiv \mathbb{R}^m.$$

b) On distingue les deux cas suivants :

- $y_1 = \dots = y_m = y$: Dans ce cas, le système linéaire admet une unique solution $\beta = y$, qui est aussi la solution du problème aux moindres carrés (le système est de rang 1).
- $\exists i \neq j, y_i \neq y_j$: dans ce cas, le système linéaire ne possède pas de solution. En revanche, il existe une infinité de solutions au sens des moindres carrés; parmi celles-ci, la solution

$$\beta^* = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y} = \frac{1}{m} \sum_{i=1}^m y_i$$

est celle de norme minimale.

Remarque : ce second cas généralise le précédent pour ce qui est de la solution au sens des moindres carrés.

c) En définissant la matrice $\mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \in \mathbb{R}^{m \times 1}$, on peut alors modéliser le problème soit comme la résolution du système linéaire

$$\begin{cases} x_1 \beta = y_1 \\ \vdots \\ x_m \beta = y_m. \end{cases}$$

soit comme celle du problème aux moindres carrés

$$\min_{\beta \in \mathbb{R}} \frac{1}{2} \sum_{i=1}^m (\beta x_i - y_i)^2 = \frac{1}{2} \|\mathbf{X}\beta - \mathbf{y}\|^2,$$

En termes de résolution, on distinguera comme en question b) le cas où le système linéaire est compatible (càd admet une solution) de celui où il ne l'est pas, pour appliquer par la suite les résultats du cours.

Note : Cette modélisation est très simpliste, mais elle a le mérite d'être explicable. D'autres modèles linéaires pourraient être proposés par les étudiants.

Solution de l'exercice 3.2: Matrices de rang 1

- a) Une décomposition de \mathbf{X} est de la forme $\mathbf{U}\Sigma\mathbf{V}^T$, avec $\mathbf{U} \in \mathbb{R}^{m \times 1}$ orthogonale, $\Sigma \in \mathbb{R}^{m \times 1}$ et $\mathbf{V} \in \mathbb{R}^{1 \times 1}$. Comme \mathbf{X} est de rang 1, cette matrice représente un vecteur $\mathbf{x} \in \mathbb{R}^m$ non nul. Soit $\mathbf{u} = \frac{\mathbf{x}}{\|\mathbf{x}\|}$ et $\mathbf{U} \in \mathbb{R}^{m \times m}$ une matrice orthogonale ayant \mathbf{u} comme première colonne. Alors, en

posant $\mathbf{V} = [1]$ (qui est bien orthogonale dans $\mathbb{R}^{1 \times 1}$!) et $\Sigma = \begin{bmatrix} \|\mathbf{x}\| \\ 0 \\ \vdots \\ 0 \end{bmatrix}$ (vecteur colonne de

taille m avec toutes les composantes nulles sauf la première), on a bien $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^T$.

- b) On cherche ici un réel β . Dans le cas général, la solution de norme minimale au problème des moindres carrés est donnée par $\hat{\beta} = \mathbf{X}^\dagger \mathbf{y}$. Comme \mathbf{X} est de rang 1 à une seule colonne, la pseudo-inverse s'écrit ici $\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$, et $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ est l'unique solution au problème des moindres carrés.
- c) \mathbf{X} est une matrice orthogonale, donc \mathbf{X} représente un vecteur \mathbf{x} de norme 1. Il s'agit donc d'une base du sous-espace engendré par \mathbf{x} . On voit alors que $\mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{X}\mathbf{X}^T \mathbf{y}^2$ est la projection de \mathbf{y} sur ce sous-espace.

NB : Si \mathbf{y} appartient à l'image de \mathbf{X} (càd \mathbf{y} est colinéaire à \mathbf{x}), la projection est égale à \mathbf{y} lui-même et le système linéaire $\mathbf{X}\beta = \mathbf{y}$ admet une unique solution ($\hat{\beta}$ est solution au sens classique). Sinon, le système linéaire n'a pas de solution (mais le problème de projection en a toujours une).

Solution de l'exercice 3.3: Moindres carrés régularisés

- a) i) Si \mathbf{X} est de rang n , alors $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{n \times n}$ est aussi de rang n , et donc $\mathbf{X}^T \mathbf{X}$ est une matrice inversible. Comme on suppose $\lambda = 0$ dans cette question, le problème (1) s'écrit $\min_{\beta \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{X}\beta - \mathbf{y}\|_2^2$, et a alors pour unique solution $\beta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$.
- ii) On fournit à titre consultatif une preuve du résultat énoncé dans la question. Si l'on considère les deux expressions $\frac{1}{2} \|\mathbf{X}\beta - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \|\beta\|_2^2$ et $\frac{1}{2} \|(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{1/2} \beta - \mathbf{z}\|_2^2$, notre but est de montrer que ces expressions sont égales à une constante (indépendante de β) près. On a:

$$\begin{aligned} \frac{1}{2} \|\mathbf{X}\beta - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \|\beta\|_2^2 &= \frac{1}{2} (\beta^T \mathbf{X}^T \mathbf{X} \beta - 2\mathbf{y}^T \mathbf{X} \beta + \mathbf{y}^T \mathbf{y}) + \frac{\lambda}{2} \beta^T \beta \\ &= \frac{1}{2} \beta^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \beta - \mathbf{y}^T \mathbf{X} \beta + \frac{1}{2} \mathbf{y}^T \mathbf{y}. \end{aligned}$$

² \mathbf{X} représente un vecteur \mathbf{x} de norme 1, donc $(\mathbf{X}^T \mathbf{X})^{-1} = \|\mathbf{x}\|^{-2} = 1$.

En parallèle, on a :

$$\begin{aligned}
 \frac{1}{2} \|(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{1/2} \boldsymbol{\beta} - \mathbf{z}\|_2^2 &= \frac{1}{2} \boldsymbol{\beta}^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{T/2} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{1/2} \boldsymbol{\beta} - \mathbf{z}^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta} + \frac{1}{2} \mathbf{z}^T \mathbf{z} \\
 &= \frac{1}{2} \boldsymbol{\beta}^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta} - [(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{1/2} \mathbf{z}]^T \mathbf{X} \boldsymbol{\beta} + \frac{1}{2} \mathbf{z}^T \mathbf{z} \\
 &= \frac{1}{2} \boldsymbol{\beta}^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta} - [\mathbf{X}^T \mathbf{y}]^T \boldsymbol{\beta} + \frac{1}{2} \mathbf{z}^T \mathbf{z} \\
 &= \frac{1}{2} \boldsymbol{\beta}^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta} - \mathbf{y}^T \mathbf{X} \boldsymbol{\beta} + \frac{1}{2} \mathbf{y}^T \mathbf{y} - \frac{1}{2} \mathbf{y}^T \mathbf{y} + \frac{1}{2} \mathbf{z}^T \mathbf{z}
 \end{aligned}$$

où l'on a utilisé le fait que $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ (et donc $(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{1/2}$, cf exo 2) est symétrique entre la première et la seconde ligne, et le lien entre \mathbf{y} et \mathbf{z} entre la seconde ligne et la troisième.

En tant que fonctions de $\boldsymbol{\beta}$, ces deux expressions ne diffèrent que par un terme constant indépendant de $\boldsymbol{\beta}$: les problèmes de minimisation associés ont donc le même ensemble de solutions (cf cours d'optimisation au second semestre).

Réponse à la question : En observant le second problème, on voit que $(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{1/2} \succ 0$ (car $\mathbf{X}^T \mathbf{X} \succ 0$) et donc que le problème admet une unique solution donnée par

$$\left[(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{T/2} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{1/2} \right]^{-1} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{T/2} \mathbf{z} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}.$$

- b) i) D'après le cours, on sait que la solution donnée par la pseudo-inverse est, parmi les solutions possibles, celle de norme minimale.
- ii) Soit $\mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T$ une décomposition en valeurs singulières de \mathbf{X} . Alors, on peut écrire une décomposition en valeurs singulières de $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ comme suit :

$$\begin{aligned}
 \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} &= \mathbf{V} \boldsymbol{\Sigma}^T \mathbf{U}^T \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T + \lambda \mathbf{I} \\
 &= \mathbf{V} \boldsymbol{\Sigma}^T \boldsymbol{\Sigma} \mathbf{V}^T + \lambda \mathbf{I} \\
 &= \mathbf{V} \boldsymbol{\Sigma}^T \boldsymbol{\Sigma} \mathbf{V}^T + \lambda \mathbf{V} \mathbf{V}^T \\
 &= \mathbf{V} (\boldsymbol{\Sigma}^T \boldsymbol{\Sigma} + \lambda \mathbf{I}) \mathbf{V}^T.
 \end{aligned}$$

Il s'agit également d'une décomposition en valeurs propres de la matrice, et les valeurs propres sont données par la matrice $\boldsymbol{\Sigma}^T \boldsymbol{\Sigma} + \lambda \mathbf{I}$. Cette matrice est diagonale à coefficients diagonaux $\boldsymbol{\Sigma}_{11}^2 + \lambda, \boldsymbol{\Sigma}_{22}^2 + \lambda, \dots, \boldsymbol{\Sigma}_{nn}^2 + \lambda$. Ces coefficients sont tous supérieurs ou égaux à $\lambda > 0$, d'où l'on peut conclure que la matrice est définie positive, et donc de rang plein.

- iii) Comme la matrice est de rang plein, il existe une unique solution au problème (1), qui est donnée par $(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$, comme pour la question a) ii).