
TD N° 2 : Optimisation sous contraintes

EXERCICE 1. Résoudre le problème suivant :

$$\begin{aligned} \min_{(x_1, x_2, x_3) \in \mathbb{R}^3} & \frac{1}{2} [(x_1 - 2)^2 + (x_2 - 2)^2 + (x_3 - 2)^2] \\ \text{s.c.} & \quad x_1 + x_2 + x_3 = 1 \end{aligned} \quad (1)$$

Correction:

$$\mathcal{L}(x, \lambda) = \frac{1}{2} [(x_1 - 2)^2 + (x_2 - 2)^2 + (x_3 - 2)^2] + \lambda(x_1 + x_2 + x_3 - 1) \quad (2)$$

Conditions du premier ordre :

$$\nabla_x \mathcal{L}(x, \lambda) = 0 \iff \begin{cases} (x_1 - 2) + \lambda = 0 \\ (x_2 - 2) + \lambda = 0 \\ (x_3 - 2) + \lambda = 0 \end{cases} \quad (3)$$

Cela assure que $x_1 = x_2 = x_3$, et donc avec la contrainte de satisfiabilité, $x_1 + x_2 + x_3 = 1$

EXERCICE 2.

On prend $y \in \mathbb{R}^n$, $X \in \mathbb{R}^{n \times p}$ avec $\text{rg}(X) = p$, $R \in \mathbb{R}^{p \times q}$ avec $\text{rg}(R) = q$ et enfin $r \in \mathbb{R}^q$. On définit les moindres carrés sous contraintes de la manière suivantes :

$$\begin{aligned} \hat{\beta}_c & \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|y - X\beta\|^2 \\ \text{s.c.} & \quad R\beta = r \end{aligned}$$

1) Montrer qu'avec $\hat{\beta}$ solution des moindres carrés (non contraints), on obtient :

$$\hat{\beta}_c = \hat{\beta} + (X^\top X)^{-1} R^\top \left[R (X^\top X)^{-1} R^\top \right]^{-1} (r - R\hat{\beta})$$

2) Montrer de plus que sous l'hypothèse gaussienne de modèle de régression linéaire classique :

$$\frac{1}{q\sigma^2} (R(\hat{\beta} - \beta^*))^\top \left[R (X^\top X)^{-1} R^\top \right]^{-1} R(\hat{\beta} - \beta^*) \sim \mathcal{F}_{n-p}^q$$

où \mathcal{F}_{n-p}^q est une loi de Fisher.

- 3) Montrer que $(R\hat{\beta} - r)^\top \left[R (X^\top X)^{-1} R^\top \right]^{-1} (R\hat{\beta} - r) = \|\hat{y} - \hat{y}_c\|^2$, et que $\|\hat{y} - \hat{y}_c\|^2 = \|y - \hat{y}_c\|^2 - \|\hat{y} - y\|^2$ où $\hat{y} = X\hat{\beta}$.
- 4) Prendre $X = [\mathbf{1}_{C_1}, \dots, \mathbf{1}_{C_K}]$, et en déduire un test de l'hypothèse " $\mu_1 = \dots = \mu_K$ " dans l'ANOVA à un facteur.

Correction:

1)

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial \beta} = -2X^\top Y + 2X^\top X\hat{\beta}_c - R^\top \hat{\lambda} = 0 \\ \frac{\partial \mathcal{L}}{\partial \lambda} = R\hat{\beta}_c - r = 0 \end{cases}$$

Ensuite :

$$\begin{aligned} -2R(X^\top X)^{-1}X^\top Y + 2R(X^\top X)^{-1}X^\top X\hat{\beta}_c - R(X^\top X)^{-1}R^\top \hat{\lambda} &= 0 \\ -2R(X^\top X)^{-1}X^\top Y + 2R\hat{\beta}_c - R(X^\top X)^{-1}R^\top \hat{\lambda} &= 0 \\ -2R(X^\top X)^{-1}X^\top Y + 2r - R(X^\top X)^{-1}R^\top \hat{\lambda} &= 0 \end{aligned}$$

Puis on en déduit :

$$\hat{\lambda} = 2 \left[R(X^\top X)^{-1}R^\top \right]^{-1} \left[r - R(X^\top X)^{-1}X^\top Y \right]$$

Ensuite, remplaçons dans la dérivée du lagrangien par rapport à β :

$$-2X^\top Y + 2X^\top X\hat{\beta}_c - 2R^\top \left[R(X^\top X)^{-1}R^\top \right]^{-1} \left[r - R(X^\top X)^{-1}X^\top Y \right] = 0$$

Enfin,

$$\begin{aligned} \hat{\beta}_c &= (X^\top X)^{-1}X^\top Y + (X^\top X)^{-1}R^\top \left[R(X^\top X)^{-1}R^\top \right]^{-1} (r - R\hat{\beta}) \\ &= \hat{\beta} + (X^\top X)^{-1}R^\top \left[R(X^\top X)^{-1}R^\top \right]^{-1} (r - R\hat{\beta}) \end{aligned}$$

- 2) $R\hat{\beta} \sim \mathcal{N}(R\beta^*, \sigma^2 R(X^\top X)^{-1}R^\top)$ avec $R(X^\top X)^{-1}R^\top \in \mathbb{R}^{q \times q}$ inversible (car X et R sont de plein rang). On rappelle que pour l'estimateur des moindres carrés $\hat{\beta}$ et $\hat{\sigma}^2$ sont indépendants ()

$$\frac{1}{\sigma^2} (R(\hat{\beta} - \beta^*))^\top \left[R(X^\top X)^{-1}R^\top \right]^{-1} R(\hat{\beta} - \beta^*) \sim \chi_q^2.$$

On termine en rappelant que $(n-p)\frac{\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p}^2$ où $\hat{\sigma}^2 = \frac{1}{n-p} \|y - X\hat{\beta}\|^2$, et que les deux χ^2 sont indépendants :

$$\hat{\beta} = (X^\top X)^{-1}X^\top Y = (X^\top X)^{-1}X^\top \left(X(X^\top X)^{-1}X^\top \right) Y = (X^\top X)^{-1}X^\top \Pi_X Y,$$

et en rappelant que par le théorème de Cochran $\Pi_X(y)$ et $(\text{Id} - \Pi_X)(y)$ sont indépendant et donc il en est de même pour les deux χ^2 .

- 3) On repart de la première question pour obtenir :

$$(X^\top X)(\hat{\beta}_c - \hat{\beta}) = R^\top \left[R(X^\top X)^{-1}R^\top \right]^{-1} (r - R\hat{\beta})$$

En multipliant $(\hat{\beta}_c - \hat{\beta})^\top$ on obtient

$$\begin{aligned} \underbrace{(\hat{\beta}_c - \hat{\beta})^\top (X^\top X)(\hat{\beta}_c - \hat{\beta})}_{=\|\hat{y}_c - \hat{y}\|^2} &= (\hat{\beta}_c - \hat{\beta})^\top R^\top \left[R(X^\top X)^{-1}R^\top \right]^{-1} (r - R\hat{\beta}) \\ &= (r - R\hat{\beta})^\top \left[R(X^\top X)^{-1}R^\top \right]^{-1} (r - R\hat{\beta}). \end{aligned}$$

Enfin on montre que

$$\begin{aligned} \|\hat{y}_c - y\|^2 &= \|\hat{y}_c - \hat{y}\|^2 + \|\hat{y} - y\|^2 + 2\langle \hat{y}_c - \hat{y}, \hat{y} - y \rangle \\ &= \|\hat{y}_c - \hat{y}\|^2 + \|\hat{y} - y\|^2 + 2\langle X\hat{\beta}_c - X\hat{\beta}, X\hat{\beta} - y \rangle \end{aligned}$$

Le dernier produit scalaire est nul car le vecteur des résidus dans les moindres carrés est orthogonal à toutes les variables explicatives.

4) Ici on a $\hat{\beta} = (\bar{y}_1, \dots, \bar{y}_K)^\top$, donc $(X^\top X)^{-1} = \text{diag}(\frac{1}{n_1}, \dots, \frac{1}{n_K})$ et donc $\hat{y} = \sum_{k=1}^K \mathbb{1}_{C_k} \bar{y}_{C_k}$. Dans ce contexte $R = \begin{bmatrix} 1 & -1 & & & \\ & 1 & -1 & & \\ & & \ddots & \ddots & \\ & & & 1 & -1 \end{bmatrix} \in \mathbb{R}^{(K-1) \times K}$ et $r = 0$ et donc $q = K - 1$. On peut aussi voir que la solution du problème

$$\begin{aligned} \hat{\beta}_c &\in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|y - X\beta\|^2 \\ \text{s.c. } &\beta_1 = \dots = \beta_K \end{aligned}$$

donne simplement $\hat{\beta}_1 = \dots = \hat{\beta}_K = \bar{y}_n$ et donc $\hat{y}_c = \bar{y}_n \mathbb{1}_n$. En rappelant que $R\beta^\star = 0$ et que $\|\hat{y} - \hat{y}_c\|^2 = \|y - \hat{y}_c\|^2 - \|\hat{y} - y\|^2$, on obtient :

$$\frac{\frac{1}{K-1} \sum_{k=1}^K (\bar{y}_{C_k} - \bar{y}_n)^2}{\frac{1}{n-K} \sum_{i=1}^n (\hat{y}_i - y_i)^2} \sim \mathcal{F}_{n-K}^{K-1}.$$

EXERCICE 3. Soit $X = \left[\frac{\mathbb{1}_{C_1}}{\sqrt{n_1}}, \dots, \frac{\mathbb{1}_{C_K}}{\sqrt{n_K}} \right]$ (avec $\mathbf{x}_k = \frac{\mathbb{1}_{C_k}}{\sqrt{n_k}}$), avec $\hat{\mu}_k = \frac{1}{n_k} \sum_{i \in C_k} y_i$ et $n_k = \#\{i \in C_k, i \in \llbracket 1, n \rrbracket\}$ avec $n_1 + \dots + n_K = n$, et on suppose que les C_k forment une partition de $\llbracket 1, n \rrbracket$.

L'estimateur Ridge (sans pénalité sur la constante) est solution de

$$(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_K)^\top = \arg \min_{\beta_0, \dots, \beta_K} \left\| \mathbf{y} - \beta_0 \mathbb{1}_n - \sum_{j=1}^K \beta_j \mathbf{x}_j \right\|^2 + \lambda \sum_{j=1}^K \beta_j^2$$

- 1) Donner la valeur de $X^\top X$ et de $X^\top y$.
- 2) Donner une formule explicite pour l'estimateur Ridge, en fonction de $y, \hat{\mu}_1, \dots, \hat{\mu}_K$, et des n_1, \dots, n_K .
- 3) Comparer avec l'estimateur Ridge classique obtenu en forçant la contrainte $\beta_0 = 0$.

Correction:

- 1) $X^\top X = \text{Id}_K$ et $X^\top \mathbf{y} = (\sqrt{n_1} \hat{\mu}_1, \dots, \sqrt{n_K} \hat{\mu}_K)^\top$.
- 2) Notons $e = (\sqrt{n_1}, \dots, \sqrt{n_K})^\top$ et vérifions que $\mathbb{1}_n^\top X = e^\top$. Partons des conditions nécessaires du premier ordre :

$$\begin{cases} \langle \mathbb{1}_n, X\tilde{\beta} + \mathbb{1}_n \hat{\beta}_0 - y \rangle = 0 \\ X^\top (X\tilde{\beta} + \mathbb{1}_n \hat{\beta}_0 - y) + \lambda \tilde{\beta} = 0 \end{cases} \iff \begin{cases} e^\top \tilde{\beta} = n\bar{y}_n - n\hat{\beta}_0 \\ (1 + \lambda)\tilde{\beta} + e\hat{\beta}_0 - (\sqrt{n_1}\hat{\mu}_1, \dots, \sqrt{n_K}\hat{\mu}_K)^\top = 0 \end{cases}$$

En multipliant la dernière équation à gauche et à droite par e on obtient :

$$(1 + \lambda)e^\top \tilde{\beta} + n\hat{\beta}_0 - \sum_{k=1}^K n_k \mu_k = (1 + \lambda)e^\top \tilde{\beta} + n\hat{\beta}_0 - n\bar{y}_n = 0$$

En y substituant l'expression de $e^\top \tilde{\beta}$ obtenue dans le système précédent, on obtient :

$$(1 + \lambda)(n\bar{y}_n - n\hat{\beta}_0) + n\hat{\beta}_0 - n\bar{y}_n = 0$$

et donc comme $\lambda > 0$,

$$\boxed{\hat{\beta}_0 = \bar{y}_n}.$$

Enfin la dernière équation du système précédent donne alors :

$$\tilde{\beta} = \frac{1}{1+\lambda} (\sqrt{n_1}(\hat{\mu}_1 - \bar{y}_n), \dots, \sqrt{n_K}(\hat{\mu}_K - \bar{y}_n)) .$$

Enfin on obtient donc comme prédicteur associé \hat{y} donné par :

$$\begin{aligned} \hat{y} &= X\tilde{\beta} + \hat{\beta}_0 \mathbf{1}_n \\ \hat{y} &= \frac{1}{1+\lambda} \sum_{k=1}^K (\hat{\mu}_k - \bar{y}_n) \mathbf{1}_{C_k} + \bar{y}_n \mathbf{1}_n \\ \hat{y} &= \frac{1}{1+\lambda} \sum_{k=1}^K (\hat{\mu}_k \mathbf{1}_{C_k} + \lambda \bar{y}_n \mathbf{1}_n)^\top . \end{aligned}$$

Ainsi cela signifie que pour un élément $i \in C_k$, le prédicteur associé est donné par :

$$\hat{y}_k = \frac{1}{1+\lambda} (\lambda \bar{y}_n + \hat{\mu}_k)$$

Interprétation : λ permet d'osciller entre le prédicteur globale (\bar{y}_n , si $\lambda = +\infty$) et le prédicteur par modalité ($\hat{\mu}_k$, si $\lambda = 0$)

3) Dans ce cas l'estimateur est défini par :

$$\arg \min_{\beta_1, \dots, \beta_K} \left\| \mathbf{y} - \sum_{j=1}^K \beta_j \mathbf{x}_j \right\|^2 + \lambda \sum_{j=1}^K \beta_j^2$$

et l'on trouve :

$$\hat{y}_k = \frac{1}{1+\lambda} \hat{\mu}_k$$

EXERCICE 4. On observe K classes C_1, \dots, C_K et n observations d'un phénomène (*e.g.*, le rendement de la variétés) sont consignées. On fait l'hypothèse que les classes C_k sont disjointes et forment une partition des observations : $\cup_{k=1}^K C_k = \llbracket 1, n \rrbracket$ et $\forall (k, k') \in \llbracket 1, K \rrbracket, C_k \cap C_{k'} = \emptyset$. Enfin on suppose que la cardinalité de chaque classe C_k est n_k , et donc que $n = \sum_{k=1}^K n_k$.

Le modèle linéaire associé peut s'écrire de la façon suivant : on définit les α_k , coefficients qui correspondent au niveau d'influence de la k^e classe

$$y = \underbrace{\begin{bmatrix} \mathbf{1}_n & \mathbf{1}_{C_1} & \dots & \mathbf{1}_{C_K} \end{bmatrix}}_X \underbrace{\begin{bmatrix} \mu \\ \alpha_1 \\ \vdots \\ \alpha_K \end{bmatrix}}_{\beta} + \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_K \end{bmatrix} .$$

Donner alors l'estimateur des moindres carrés sous contraintes associé :

$$\begin{aligned} \min_{\beta \in \mathbb{R}^{K+1}} \frac{1}{2} \|y - X\beta\|^2 , \\ \text{t.q. } \beta = (\mu, \alpha_1, \dots, \alpha_K)^\top \text{ et } \sum_{k=1}^K c_k \alpha_k = 0 , \end{aligned}$$

où le vecteur $c = (c_1, \dots, c_K)^\top \in \mathbb{R}^K$ est un vecteur encodant les contraintes choisies tel que $\sum_{k=1}^K c_k \neq 0$.

EXERCICE 5. Soit $X = \left[\frac{\mathbf{1}_{C_1}}{\sqrt{n_1}}, \dots, \frac{\mathbf{1}_{C_K}}{\sqrt{n_K}} \right]$ (avec $\mathbf{x}_k = \frac{\mathbf{1}_{C_k}}{\sqrt{n_k}}$), avec $\hat{\mu}_k = \frac{1}{n_k} \sum_{i \in C_k} y_i$ et $n_k = \#\{i \in C_k, i \in \llbracket 1, n \rrbracket\}$ avec $n_1 + \dots + n_K = n$.

L'estimateur Lasso (sans pénalité sur la constante) est solution de

$$(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_K)^\top = \arg \min_{\beta_0, \dots, \beta_K} \frac{1}{2} \left\| \mathbf{y} - \beta_0 \mathbf{1}_n - \sum_{j=1}^K \beta_j \mathbf{x}_j \right\|_2^2 + \lambda \sum_{j=1}^K |\beta_j|$$

- 1) Donner la valeur de $X^\top X$ et de $X^\top y$
- 2) Donner une formule explicite pour l'estimateur Lasso, en fonction de $y, \hat{\mu}_1, \dots, \hat{\mu}_K$, et des n_1, \dots, n_K .