

Exercice 1 :

Exercice Class : Linear Regression

1) Soit $y \in \mathbb{R}^m$ alors $\text{var}_m(y) = \frac{1}{m} \sum_{i=1}^m (y_i - \bar{y}_m)^2$ avec $\bar{y}_m = \frac{1}{m} \sum_{i=1}^m y_i$

$$\text{var}_m(y) = 0 \Leftrightarrow \frac{1}{m} \sum_{i=1}^m (y_i - \bar{y}_m)^2 = 0 \Leftrightarrow \sum_{i=1}^m \underbrace{(y_i - \bar{y}_m)^2}_{\geq 0} = 0$$

$$\Leftrightarrow \forall i=1, \dots, m \quad y_i = \bar{y}_m$$

$$\boxed{\text{var}_m(y) = 0 \Leftrightarrow y \in \text{Span}(\mathbb{1}_m) \Leftrightarrow \exists c \in \mathbb{R} / \forall i, y_i = c}$$

2) Montrons que $\text{Ker}(X) = \text{Ker}(X^T X)$ avec $X \in \mathbb{R}^{m \times p}$

[C]. Soit $u \in \text{Ker}(X)$ alors $Xu = 0$ donc $X^T Xu = 0$ et $u \in \text{Ker}(X^T X)$

[D]. Soit $u \in \text{Ker}(X^T X)$ alors $X^T Xu = 0$ puis $u^T X^T Xu = 0$, i.e. $\|Xu\|^2 = 0$
donc $Xu = 0$ et $u \in \text{Ker}(X)$, ainsi $\boxed{\text{Ker}(X) = \text{Ker}(X^T X)}$

3) $\mathbb{1}_m = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{R}^{m \times 1}$, $\tilde{X} \in \mathbb{R}^{m \times p}$, $X = (\mathbb{1}_m, \tilde{X}) \in \mathbb{R}^{m \times (p+1)}$

$$\tilde{X} = \begin{pmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_m^T \end{pmatrix} \in \mathbb{R}^{m \times p} \quad \hat{\mu}_m = \frac{1}{m} \sum_{i=1}^m \tilde{x}_i$$

$$(X^T X) \text{ non inversible} \Leftrightarrow \text{Ker}(X^T X) \neq \{0\} \Leftrightarrow \text{Ker}(X) \neq \{0\}$$

$$\Leftrightarrow \exists u \in \mathbb{R}^{p+1} \setminus \{0\} / Xu = 0$$

$$\Leftrightarrow \exists u_0 \in \mathbb{R}, \tilde{u} \in \mathbb{R}^p \setminus \{0\} / u_0 \mathbb{1}_m + \tilde{X} \tilde{u} = 0$$

Posez $y = \tilde{X} \tilde{u} \in \mathbb{R}^m$

$$(X^T X) \text{ non-inversible} \Leftrightarrow \exists u_0 \in \mathbb{R}, \tilde{u} \in \mathbb{R}^p \setminus \{0\} / y = \tilde{X} \tilde{u} = -u_0 \mathbb{1}_m$$

$$\Leftrightarrow \exists \tilde{u} \in \mathbb{R}^p \setminus \{0\} / y \in \text{Span}(\mathbb{1}_m), \text{ i.e. } y = \bar{y} \mathbb{1}_m$$

Or $y = \tilde{X} \tilde{u}$ donc $\bar{y} = \bar{X}^T \tilde{u} = \hat{\mu}_m \tilde{u}$

$$(X^T X) \text{ non-inversible} \Leftrightarrow \exists \tilde{u} \in \mathbb{R}^p \setminus \{0\} / (\tilde{X} - \hat{\mu}_m) \tilde{u} = 0$$

$$\Leftrightarrow \text{Ker}(\tilde{X}^c) \neq \{0\}$$

$$\Leftrightarrow \text{Ker}(\tilde{X}^{cT} \tilde{X}^c) \neq \{0\}$$

$$= \tilde{X}^c = \begin{pmatrix} \tilde{x}_1 - \hat{\mu}_m \\ \vdots \\ \tilde{x}_m - \hat{\mu}_m \end{pmatrix}$$

$$\boxed{(X^T X) \notin GL_{p+1}(\mathbb{R}) \Leftrightarrow (\tilde{X}^{cT} \tilde{X}^c) \text{ non-inversible}}$$

$$\text{et } \tilde{X}^{cT} \tilde{X}^c = \sum_{i=1}^m (\tilde{x}_i - \hat{\mu}_m)(\tilde{x}_i - \hat{\mu}_m)^T$$

4) Sur les variables centrées $\tilde{X} \in \mathbb{R}^{m \times p}$, $Y^c \in \mathbb{R}^m$,
 $\hat{\theta}_m \in \arg \min_{\theta \in \mathbb{R}^p} \|Y^c - \tilde{X}\theta\|_2^2$

L'estimateur est donné par : (on a supposé $\text{cov}_m(\tilde{X})$ inversible)

$$\hat{\theta}_m = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T Y^c = \frac{1}{n} \text{cov}_m(\tilde{X})^{-1} \tilde{X}^T Y^c$$

Exercice 2:

$\tilde{X} \in \mathbb{R}^{m \times p}$ de rang plein $\tilde{X} = (\tilde{X}_1 | \dots | \tilde{X}_p)$, $\tilde{X}_k \in \mathbb{R}^m \forall k=1, \dots, p$

1) On change d'échelle pour la colonne k : $\tilde{X}_k \leftarrow b \tilde{X}_k$ avec $b > 0$

Posons $D = \text{Diag}(1, \dots, 1, \underbrace{b}_{\text{position } (k+1)}, 1, \dots, 1)$, $D \in \mathbb{R}^{(p+1) \times (p+1)}$

$$XD = (1_n, \tilde{X})D = \left(\begin{array}{c|c} 1 & \tilde{X}_1 \dots \tilde{X}_p \end{array} \right) D = \left(\begin{array}{c|c} 1 & \tilde{X}_1 \times 1 \dots \tilde{X}_k \times b \dots \tilde{X}_p \times 1 \end{array} \right)$$

$$XD = \left(\begin{array}{c|c} 1 & \tilde{X}_1 \dots \tilde{X}_k b \dots \tilde{X}_p \end{array} \right) = (1_n, \tilde{X}_1, \dots, \tilde{X}_k b, \dots, \tilde{X}_p) = X_b$$

$$D = \text{Diag}(1, \dots, 1, b, 1, \dots, 1), \quad XD = X_b$$

2) Par définition:

$$\hat{\theta}_{b,m} \in \arg \min_{\theta \in \mathbb{R}^{p+1}} \{ \|Y - X_b \theta\|_2^2 = \|Y - XD \theta\|_2^2 \}$$

d'après les équations normales (projection, Hilbert)

$$(\tilde{X}_b^T X_b) \hat{\theta}_{b,m} = \tilde{X}_b^T Y$$

i.e.

$$D(\tilde{X}^T X) D \hat{\theta}_{b,m} = D \tilde{X}^T Y$$

donc

$$(\tilde{X}^T X) D \hat{\theta}_{b,m} = \tilde{X}^T Y$$

ainsi,

$$\hat{\theta}_m = D \hat{\theta}_{b,m}$$

3) $\hat{\theta}_{b,m} = \tilde{D}^T \hat{\theta}_m$ donc $\text{var}(\hat{\theta}_{b,m}) = \tilde{D}^T \text{var}(\hat{\theta}_m) \tilde{D}$

model fixed design: $Y = X\theta^* + \epsilon$, $\text{var}(\epsilon) = \sigma^2 I_m$, $\text{var}(\hat{\theta}_m) = (\tilde{X}^T \tilde{X})^{-1} \sigma^2$

$$\text{var}(\hat{\theta}_{b,m}) = \sigma^2 \tilde{D}^T (\tilde{X}^T \tilde{X})^{-1} \tilde{D}$$

4) On regarde la valeur prédite par le modèle OLS

Pour le modèle initial (X, Y) , la fonction de prédiction est

$$g: \mathbb{R}^p \rightarrow \mathbb{R}, \quad g(x) = x^T \hat{\theta}_m$$

Pour le modèle modifié, (X_b, Y) , la fonction de prédiction est

$$g_b: \mathbb{R}^p \rightarrow \mathbb{R}, \quad g_b(x) = x_b^T \hat{\theta}_{b,m} = x^T (D D^T) \hat{\theta}_m = x^T \hat{\theta}_m$$

Pour la prédiction, on a $\boxed{g_b(x) = g(x)}$

5) On regarde la valeur prédite par le modèle RIDGE

Pour le modèle modifié (X_b, Y) la fonction de prédiction est

$$g_b^{(rdg)}(x) = x_b^T \hat{\theta}_{b,m}^{(rdg)} = x^T D \hat{\theta}_m^{(rdg)}$$

avec $\hat{\theta}_{b,m}^{(rdg)} \in \arg \min_{\theta \in \mathbb{R}^p} \{ \|Y - X_b \theta\|_2^2 + \lambda \|\theta\|_2^2 \} = (X_b^T X_b + \lambda I)^{-1} X_b^T Y$

donc $\boxed{g_b^{(rdg)}(x) \neq g^{(rdg)}(x)}$ $\hat{\theta}_{b,m}^{(rdg)} \neq D^T \hat{\theta}_m^{(rdg)}$

Il n'y a plus cette propriété d'invariance, c'est pourquoi on doit renormaliser les variables avec l'estimateur RIDGE

Exercice 3:

1) La fonction de prédiction OLS est : $\boxed{\hat{p}(x) = x^T \hat{\theta}_m}$

2) $\hat{\theta}_m = (X^T X)^{-1} X^T Y$ avec $Y = X \theta^* + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$

ainsi $Y \sim \mathcal{N}(X \theta^*, \sigma^2 I_m)$

puis $\hat{\theta}_m$ est aussi gaussien avec $E[\hat{\theta}_m] = (X^T X)^{-1} X^T (X \theta^*) = \theta^*$
et variance $\sigma^2 (X^T X)^{-1}$, $\hat{\theta}_m \sim \mathcal{N}(\theta^*, \sigma^2 (X^T X)^{-1})$

$$\boxed{\hat{p}(x) \sim \mathcal{N}(x^T \theta^*, x^T (X^T X)^{-1} x \sigma^2)}$$

$$3) \sum_{i=1}^m (y_i - x_i^T \hat{\theta}_m)^2 / \sigma^2 \sim \chi_{m-(p+1)}^2 \quad ; \quad \hat{\sigma}_m^2 (m-(p+1)) = \sum_{i=1}^m (y_i - x_i^T \hat{\theta}_m)^2$$

$$\frac{\hat{p}(x) - p(x)}{\hat{\sigma}_m \sqrt{x^T (X^T X)^{-1} x}} = \frac{\hat{p}(x) - p(x)}{\sqrt{v(x)}} \times \frac{1}{\sqrt{\hat{\sigma}_m^2 / \sigma^2}} \sim T_{m-(p+1)}$$

$\underbrace{\quad}_{\mathcal{N}(0,1)} \quad = \quad \underbrace{\quad}_{= 2 \sqrt{\frac{\chi_{m-(p+1)}^2}{m-(p+1)}}}$

4) Définir $\text{var}(x) = x^T (X^T X)^{-1} x$

$$\frac{y - \hat{p}(x)}{\hat{\sigma}_m \sqrt{1 + \text{var}(x)}} = \frac{(y - p(x)) + (p(x) - \hat{p}(x))}{\hat{\sigma}_m \sqrt{1 + \text{var}(x)}} = \frac{\varepsilon + (p(x) - \hat{p}(x))}{\hat{\sigma}_m \sqrt{1 + \text{var}(x)}}$$

Au numérateur, somme de deux gaussiennes indépendantes, donc gaussienne moyenne nulle et variance $\sigma^2(1 + \text{var}(x))$, c'est la même écriture que Q3.

5) Pour $p(x)$: $\left[\hat{p}(x) - \hat{\sigma}_m \sqrt{\text{var}(x)} t_{1-\alpha/2}; \hat{p}(x) + \hat{\sigma}_m \sqrt{\text{var}(x)} t_{1-\alpha/2} \right]$

avec $t_{1-\alpha/2}$ quantile $(1-\frac{\alpha}{2})$ de la Student $(n-p+1)$.

• Pour Y , on remplace $\text{var}(x)$ par $(1 + \text{var}(x))$

Exercice 4:

1) $\hat{\theta}_m^{\text{rdg}} \in \underset{\theta \in \mathbb{R}^p}{\text{argmin}} \{ \|Y - X\theta\|_2^2 + \lambda \|\theta\|_2^2 \} = \underset{\theta \in \mathbb{R}^p}{\text{argmin}} f(\theta)$

$\theta \mapsto f(\theta)$ strictement convexe, coercive, existence et unicité du minimum.

Condition du premier ordre ($\nabla f(\theta) = 0$) fournit: $\hat{\theta}_m^{\text{rdg}} = (X^T X + \lambda I_p)^{-1} X^T Y$

2) $Y = X\theta^* + \varepsilon$ avec $E[Y] = X\theta^*$ donc

$$E[\hat{\theta}_m^{\text{rdg}}] = (X^T X + \lambda I_p)^{-1} X^T (X\theta^*) = \theta^* - \lambda (X^T X + \lambda I_p)^{-1} \theta^*$$

$$\text{Var}(\hat{\theta}_m^{\text{rdg}}) = \sigma^2 (X^T X + \lambda I_p)^{-1} X^T X (X^T X + \lambda I_p)^{-1}$$

3) L'estimateur ridge s'écrit: $\hat{\theta}_m^{\text{rdg}} = AY$ avec Y gaussien donc $\hat{\theta}_m^{\text{rdg}}$ est gaussien avec moyenne et variance ci-dessus.

4) $\frac{\hat{\theta}_{m,k}^{\text{rdg}} - \theta_k^*}{\sqrt{\text{var}(\hat{\theta}_{m,k}^{\text{rdg}})}} = \frac{\hat{\theta}_{m,k}^{\text{rdg}} - \theta_k^*}{\sqrt{v_k}} \sim N(0, 1)$

avec $v_k = e_k^T (X^T X + \lambda I_p)^{-2} (X^T X) e_k$
 $e_k^T = (0, \dots, 0, 1, 0, \dots, 0)$

L'intervalle de confiance de niveau α est alors

$$\hat{\theta}_{m,k}^{\text{rdg}} \pm \sqrt{v_k} \Phi^{-1}(1 - \alpha/2)$$