

① Leçon n° 2: Régression Linéaire / Régression Ridge

Contexte: $(x_1, y_1), \dots, (x_n, y_n) \in (\mathcal{X} \times \mathcal{Y})^n$

Minimisation du risque $R(\theta) = \mathbb{E}_{(x,y)} \left[(y - f_\theta(x))^2 \right]$, $\theta \in \mathbb{R}^d$

Estimateur de minimisation du risque empirique

$$\hat{\theta} \in \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - f_\theta(x_i))^2 =: \hat{R}(\theta)$$

Modèles Linéaires $f_\theta(x) = \underbrace{\varphi(x)}^{\text{feature map}} \theta$; $\theta \in \mathbb{R}^d$

Exemples:

- $\varphi(x) = x$: Régression Linéaire

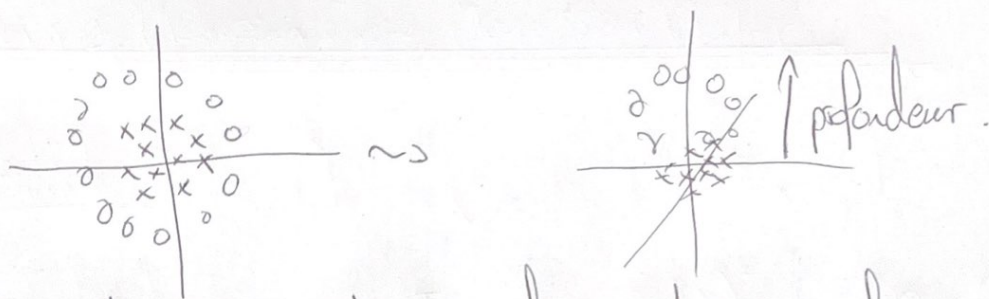
- $\varphi(x) = \begin{pmatrix} x \\ 1 \end{pmatrix}$: Régression affine (Souvent simplement appelé Régression Linéaire)

- $\varphi(x) = (\text{monômes en } x \text{ jusqu'au degré } q)$:

Régression polynomiale de degré q

- Avec de bonnes features, il est possible de rendre séparable linéairement des formes complexes (pour plus tard en classification)

②



- Avec des noyaux: peut se généraliser en dimension infinie

Notation Matricielle: En notant $y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$ et $\Phi = \begin{pmatrix} \phi(x_1) \\ \vdots \\ \phi(x_n) \end{pmatrix}$,

$$\hat{M}(\theta) = \frac{1}{n} \|y - \Phi\theta\|_2^2$$

I. Moindres Carrés:

Question: Comment trouver $\hat{\theta} \in \underset{\text{argmin}}{\hat{M}}(\theta)$?

Hypothèse 1: Φ a ses colonnes indépendantes

i.e. Il n'y a pas de dépendance linéaire entre les features du jeu de données.

En pratique, c'est le presque sûrement le cas (sous certaines hypothèses sur ϕ et sur la distribution) dès lors que n est assez grand.

③

Définition - Proposition :

Sous l'hypothèse 1, le ~~problème~~ $\hat{M}(\cdot)$ admet un unique minimum sur \mathbb{R}^d appelé estimateur des moindres carrés ordinaire.

Cet estimateur satisfait l'équation normale :

$$\Phi^T \Phi \hat{\theta} = \Phi^T y$$

Il a donc pour expression

$$\hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T y$$

preuve : Sous H1, $\hat{M}(\cdot)$ est coercive, et elle est continue et même C^1 .

$\hat{M}(\cdot)$ admet donc un minimiseur global qui est solution de l'équation $\nabla \hat{M}(\hat{\theta}) = 0$.

i.e. $\frac{2}{n} \Phi^T (\Phi \hat{\theta} - y) = 0$

ou encore $\Phi^T \Phi \hat{\theta} = \Phi^T y$.

Sous H1, cette équation admet une seule solution

$$\hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T y$$

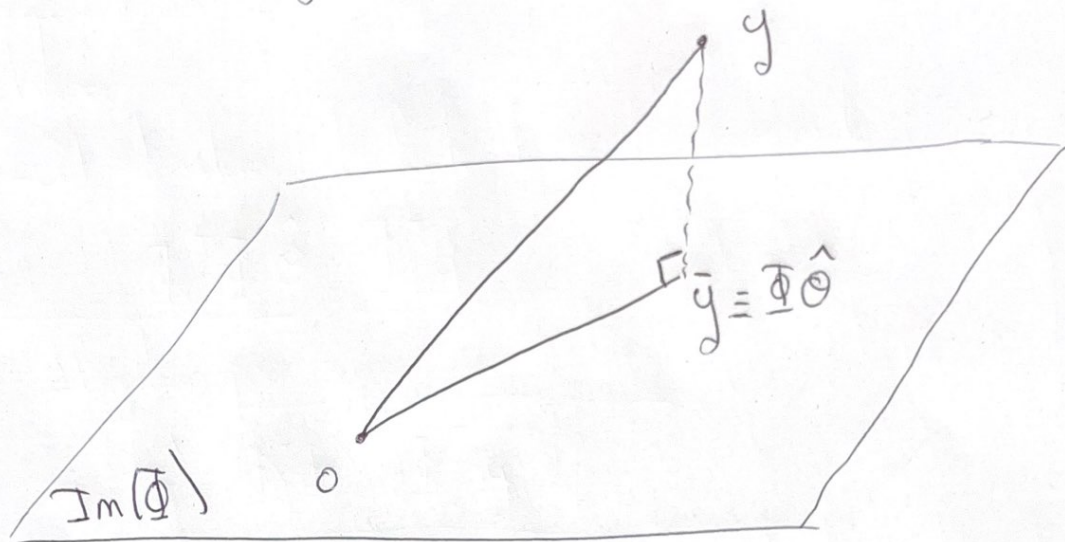
qui est à la fois un minimum local et global.

□

②

Interprétation géométrique :

- Nous avons l'orthogonalité suivante :



preuve : $\hat{\theta} \in \arg \min_{\theta} \frac{1}{n} \|y - \Phi \theta\|_2^2$

iff $\bar{y} \in \arg \min_{s \in \text{Im}(\Phi)} \frac{1}{n} \|y - s\|_2^2$
 $\underbrace{s \in \text{Im}(\Phi)}_{\in \text{Im}(\Phi)} \quad \underbrace{s \in \text{Im}(\Phi)}$

iff $\bar{y} \in \arg \min_{s \in \text{Im}(\Phi)} \frac{1}{n} \|(y - p(y)) + (p(y) - s)\|_2^2$

où $p(y)$ est la projection orthogonale de y sur $\text{Im}(\Phi)$.

iff $\bar{y} \in \arg \min_{s \in \text{Im}(\Phi)} \left\{ \frac{1}{n} \left(\|y - p(y)\|_2^2 + \|p(y) - s\|_2^2 \right) \right\}$

d'après Pythagore.

iff $\bar{y} = p(y)$.

□

⑤

II. Analyse statistique - design fixe

• Simplification: Dans les couples (x, y) , seul le y est aléatoire. Le x peut être vu comme déterministe.

Modèle:
• $\forall i, y_i = \Phi(x_i)^T \theta + \varepsilon_i$
• (ε_i) sont indépendantes
• $\forall i, \mathbb{E}(\varepsilon_i) = 0, \mathbb{E}(\varepsilon_i^2) = \sigma^2$

• Risque: $R(\theta) \equiv \frac{1}{n} \|y - \Phi \theta\|_2^2$

Proposition: $R^* = \sigma^2$. De plus, pour tout $\hat{\theta}$,

$R(\hat{\theta}) - R^* = \|\hat{\theta} - \theta_*\|_{\hat{\Sigma}}^2$ où $\hat{\Sigma} = \frac{1}{n} \Phi^T \Phi$.
Alors, $\mathbb{E}(R(\hat{\theta})) - R^* = \underbrace{\|\mathbb{E}(\hat{\theta}) - \theta_*\|_{\hat{\Sigma}}^2}_{\text{Biais}^2} + \underbrace{\mathbb{E}\left(\|\hat{\theta} - \mathbb{E}(\hat{\theta})\|_{\hat{\Sigma}}^2\right)}_{\text{Variance}}$

preuve:

$$\begin{aligned} R(\hat{\theta}) &= \frac{1}{n} \mathbb{E}_y \left(\|y - \Phi \hat{\theta}\|_2^2 \right) \\ &= \frac{1}{n} \mathbb{E}_\varepsilon \left(\frac{1}{n} \|\Phi \theta_* + \varepsilon - \Phi \hat{\theta}\|_2^2 \right) \\ &= \frac{1}{n} \mathbb{E}_\varepsilon \left(\|\Phi(\theta_* - \hat{\theta})\|_2^2 + \|\varepsilon\|_2^2 + 2(\Phi(\theta_* - \hat{\theta}))^T \varepsilon \right) \\ &= \frac{1}{n} \|\hat{\theta} - \theta_*\|_{\hat{\Sigma}}^2 + \sigma^2 + 0 \\ &= \underbrace{R(\theta_*)}_{R^*} = R^* \end{aligned}$$

⑥ De plus,

$$\begin{aligned} \bullet \quad E_{\theta}(\mathcal{M}(\hat{\theta})) - \mathcal{M}^* &= E_{\theta} \left(\|\hat{\theta} - E(\hat{\theta}) + E(\hat{\theta}) - \theta_*\|_{\hat{\Sigma}}^2 \right) \\ &= E_{\theta} \left(\|\hat{\theta} - E(\hat{\theta})\|_{\hat{\Sigma}}^2 \right) + \underbrace{2 E_{\theta} \left((\hat{\theta} - E(\hat{\theta}))^T \hat{\Sigma} (E(\hat{\theta}) - \theta_*) \right)}_{=0} \\ &\quad + E_{\theta} \left(\|E(\hat{\theta}) - \theta_*\|_{\hat{\Sigma}}^2 \right) \end{aligned}$$

• Avec l'estimateur des moindres carrés ordinaire

$$\text{Ici: } \hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T y = \hat{\Sigma}^{-1} \left(\frac{1}{n} \Phi^T y \right)$$

Analyse du biais:

$$E(\hat{\theta}) = E \left(\hat{\Sigma}^{-1} \left(\frac{1}{n} \Phi^T y \right) \right) = \hat{\Sigma}^{-1} \left(\frac{1}{n} \Phi^T E(y) \right)$$

• car: $y = \Phi^T \theta_* + \varepsilon$, donc $E(y) = \Phi^T \theta_*$ et

$$E(\hat{\theta}) = \underbrace{\hat{\Sigma}^{-1}}_I \hat{\Sigma} \theta_* = \theta_*$$

$\hat{\theta}$ est non biaisé

⑦

Analyse de la variance

$$\begin{aligned} \text{Var}(\hat{\theta}) &= E\left((\hat{\theta} - \theta_*) (\hat{\theta} - \theta_*)^T\right) \\ &= E\left((\Phi^T \Phi)^{-1} \Phi^T \varepsilon \varepsilon^T \Phi (\Phi^T \Phi)^{-1}\right) \\ &= (\Phi^T \Phi)^{-1} \Phi^T \underbrace{E(\varepsilon \varepsilon^T)}_{= \sigma^2 I} \Phi (\Phi^T \Phi)^{-1} \end{aligned}$$

$$= \sigma^2 (\Phi^T \Phi)^{-1}$$

$$= \frac{\sigma^2}{n} \hat{\Sigma}^{-1}$$

Risque de l'estimateur des moindres carrés :

Nous pouvons donc appliquer la proposition précédente et obtenir :

$$\begin{aligned} E(R(\hat{\theta})) - R^* &= E\left(\|\hat{\theta} - \theta_*\|_{\varepsilon}^2\right) \\ &= E\left((\hat{\theta} - \theta_*)^T \hat{\Sigma} (\hat{\theta} - \theta_*)\right) \\ &= E\left(\text{tr}\left(\begin{array}{c} \phantom{\hat{\Sigma}} \\ \phantom{\hat{\Sigma}} \end{array}\right)\right) \\ &= E\left(\text{tr}\left(\hat{\Sigma} (\hat{\theta} - \theta_*) (\hat{\theta} - \theta_*)^T\right)\right) \\ &= \text{tr}\left(\hat{\Sigma} E\left((\hat{\theta} - \theta_*) (\hat{\theta} - \theta_*)^T\right)\right) \\ &= \frac{\sigma^2}{n} \text{tr}(I_d) = \frac{\sigma^2 d}{n} \end{aligned}$$

③

$$\boxed{\mathbb{E}(n(\hat{\sigma}^2)) - M = \frac{\sigma^2 d}{n}}$$

III Régression de Ridge - Design Fine

Problèmes de la régression Fine et grande dimension ($d \gg n$)

- Overfitting: l'estimateur apprend le jeu d'apprentissage par cœur
- Solutions complexes car l'équation normale n'admet plus une unique solution.

Solutions:

- Réduction de dimension
- Régularisation (Régression Ridge, Régression Lasso).

Régression Ridge: Pour $\lambda > 0$,

$$\hat{\theta}_\lambda \in \operatorname{argmin}_{\theta \in \mathbb{R}^d} \frac{1}{n} \|y - \Phi \theta\|_2^2 + \lambda \|\theta\|_2^2$$

1) Expression de la solution:

Comme pour l'estimateur des moindres carrés ordinaires, il est facile de trouver une condition nécessaire sur $\hat{\theta}_\lambda$ avec de l'analyse de base. Comme $\theta \mapsto \frac{1}{n} \|y - \Phi \theta\|_2^2 + \lambda \|\theta\|_2^2$ est coercive et C¹,

$\hat{\theta}_\lambda$ existe et satisfait $\nabla \hat{M}_\lambda(\hat{\theta}_\lambda) = 0$,

i.e.
$$\frac{2}{n} (\Phi^T \Phi \hat{\theta}_\lambda - \Phi^T y) + 2\lambda \hat{\theta}_\lambda = 0.$$

g. Comme $\lambda > 0$, cette équation admet une unique solution

$$\hat{\theta}_\lambda = \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} \Phi^T y$$

2) Analyse statistique

Analyse du biais

$$\begin{aligned} E(\hat{\theta}_\lambda) &= \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} \Phi^T E(y) \\ &= \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} \Phi^T \Phi \theta_* \\ &= \cancel{\frac{1}{n}} (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} \theta_* \\ &= (\hat{\Sigma} + \lambda I)^{-1} [(\hat{\Sigma} + \lambda I) - \lambda I] \theta_* \\ &= \theta_* - \lambda (\hat{\Sigma} + \lambda I)^{-1} \theta_* \end{aligned}$$

De plus, remarquons que $\hat{\Sigma}$ et $(\hat{\Sigma} + \lambda I)^{-1}$ commutent
(~~Plus~~ $\hat{\Sigma}$ est PSD)

$$\text{Avec } \text{Biais}^2 = \|E(\hat{\theta}_\lambda) - \theta_*\|_{\hat{\Sigma}}^2 = \lambda^2 \theta_*^T (\hat{\Sigma} + \lambda I)^{-2} \hat{\Sigma} \theta_*$$

(10)

Analyse de la variance:

$$\begin{aligned}
 \text{Variance} &= E(\|\hat{\sigma}_\lambda - E(\hat{\sigma}_\lambda)\|_{\hat{\Sigma}}^2) \\
 &= E\left(\left\|\frac{1}{n}(\hat{\Sigma} + \lambda I)^{-1} \Phi^T \varepsilon\right\|_{\hat{\Sigma}}^2\right) \\
 &= E\left(\frac{1}{n^2} \text{tr}\left(\varepsilon^T \Phi (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda I)^{-1} \Phi^T \varepsilon\right)\right) \\
 &= E\left(\frac{1}{n^2} \text{tr}\left(\varepsilon \varepsilon^T \Phi (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda I)^{-1} \Phi^T\right)\right) \\
 &= \frac{\sigma^2}{n} \text{tr}\left(\hat{\Sigma} (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda I)^{-1}\right) \\
 &= \frac{\sigma^2}{n} \text{tr}\left(\hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2}\right)
 \end{aligned}$$

Analyse du risque

Par la décomposition Bias-Variance, nous avons

$$\begin{aligned}
 E(M(\hat{\sigma}_\lambda)) - M^* &= \lambda^2 \theta_v^T (\hat{\Sigma} + \lambda I)^{-2} \hat{\Sigma} \theta_v \\
 &+ \underbrace{\frac{\sigma^2}{n} \text{tr}\left(\hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2}\right)}_{\text{degrees of freedom}}
 \end{aligned}$$

(11)

Simplification de la borne supérieure

Remarquons que μ est une valeur propre de $\hat{\Sigma}$ ssi

$\frac{\lambda\mu}{(\mu+\lambda)^2}$ est une valeur propre de $(\hat{\Sigma} + \lambda I)^{-2} \lambda \hat{\Sigma}$.

Comme pour tous λ, μ nous avons $2\lambda\mu \leq (\lambda+\mu)^2$,

toutes les valeurs propres de $(\hat{\Sigma} + \lambda I)^{-2} \lambda \hat{\Sigma}$ sont inférieures à $\frac{1}{2}$.

Alors $\text{Biais}^2 \leq \frac{\lambda}{2} \|\theta_0\|^2$

et $\text{Variance} \leq \frac{\sigma^2 \text{tr}(\hat{\Sigma}^{-1})}{2\lambda n}$

Donc

$$\mathbb{E}(M(\hat{\theta}_{\lambda})) - M^0 \leq \frac{\lambda}{2} \|\theta_0\|_2^2 + \frac{\sigma^2 \text{tr}(\hat{\Sigma}^{-1})}{2\lambda n}$$

Minimiser cette expression en λ donne

$$\lambda_{\text{opt}} = \frac{\sigma \text{tr}(\hat{\Sigma}^{-1})^{1/2}}{\|\theta_0\|_2 \sqrt{n}} \quad \text{et}$$

$$\mathbb{E}(M(\hat{\theta}_{\lambda_{\text{opt}}})) - M^0 \leq \frac{\sigma \text{tr}(\hat{\Sigma}^{-1})^{1/2} \|\theta_0\|_2}{\sqrt{n}}$$
