

IMA205 - TP1 - Theoretical Questions

Théo ROUVET

5 février 2020

J'ai parfois noté p la taille de l'espace des features au lieu de d . On peut donc considérer dans ce qui suit que $d = p$.

1 Ordinary Least Squares

1.1 Espérance de $\tilde{\beta}$

$$\mathbb{E}(\tilde{\beta}) = \mathbb{E}(Cy) = \mathbb{E}((H + D)y) = \mathbb{E}(\beta^*) + \mathbb{E}(Dy) = \beta + \mathbb{E}(Dy)$$

On peut considérer un **design du type fixe** $y = x\beta + \epsilon$, tel que ϵ est un bruit blanc additif (d'espérance nulle). En injectant cette expression, on obtient :

$$\mathbb{E}(\tilde{\beta}) = \beta + \mathbb{E}(D(x\beta + \epsilon)) = \beta + \mathbb{E}(D)x\beta + \mathbb{E}(D\epsilon)$$

On peut supposer de plus que D est **déterministe** pour pouvoir continuer le calcul. On obtient par la suite :

$$\mathbb{E}(\tilde{\beta}) = \beta + Dx\beta$$

Ainsi, compte-tenu des hypothèses précédentes, $\tilde{\beta}$ est **non biaisé si et seulement si** $Dx\beta = 0$.

1.2 Variance de $\tilde{\beta}$

Calculons à présent la variance de $\tilde{\beta}$ pour montrer que l'OLS est l'estimateur non biaisé qui donne la plus petite variance.

En supposant que β^* et Dy sont **indépendants**, on écrit :

$$\begin{aligned} \text{Var}(\tilde{\beta}) &= \text{Var}(\beta^* + Dy) \\ &= \text{Var}(\beta^*) + D\text{Var}(y)D^T \\ &= (x^T x)^{-1} x^T \text{Var}(y) x (x^T x)^{-1} + D\text{Var}(y)D^T \end{aligned}$$

Or, on peut écrire que $\text{Var}(y) = \text{Var}(x\beta + \epsilon) = \text{Var}(\epsilon) = \sigma^2 I_p$ en **reprenant l'hypothèse du bruit blanc** précédente.

On a ensuite $\text{Var}(\tilde{\beta}) = \sigma^2 (x^T x)^{-1} I_p + \sigma^2 D D^T$ i.e. $\text{Var}(\tilde{\beta}) = \text{Var}(\beta^*) + \sigma^2 D D^T$

Or, $D D^T$ est une matrice symétrique et positive, et $\sigma^2 \geq 0$, **donc** $\text{Var}(\tilde{\beta}) \geq \text{Var}(\beta^*)$.

L'inégalité est stricte si et seulement si $DD^T \neq 0$ et $\sigma \neq 0$.

On obtient bien que l'OLS est l'estimateur non biaisé qui donne la plus petite variance.

2 Ridge regression

2.1 Question 1

On sait que l'estimateur Ridge est tel que $\beta_{\text{ridge}}^* = (x^T x + \lambda I_p)^{-1} x^T y$, donc par passage à l'espérance, $\mathbb{E}(\beta_{\text{ridge}}^*) = \mathbb{E}((x^T x + \lambda I_p)^{-1} x^T y) = (x^T x + \lambda I_p)^{-1} x^T x \beta$ en reprenant le design d'OLS, avec x déterministe et $y = x\beta + \epsilon$ avec ϵ centré.

Si $\lambda = 0$, on trouve une OLS et $\mathbb{E}(\beta_{\text{ridge}}^*) = \beta$.

Si $\lambda > 0$, il s'agit bien d'un Ridge et on a $(x^T x + \lambda I_p)^{-1} x^T x \neq I_n$ ce qui entraîne $\mathbb{E}(\beta_{\text{ridge}}^*) \neq \beta$ i.e. β_{ridge}^* est biaisé.

2.2 Question 2

On résout en utilisant la décomposition SVD :

$$\begin{aligned}\beta_{\text{ridge}}^* &= (x^T x + \lambda I_p)^{-1} x^T y \\ &= (VDU^T UDV^T + \lambda I_p)^{-1} x^T y \\ &= (VD^2 V^T + \lambda VV^T)^{-1} x^T y \\ &= V(D^2 + \lambda I_p)^{-1} V^T VDU^T y \\ \beta_{\text{ridge}}^* &= V(D^2 + \lambda I_p)^{-1} DU^T y\end{aligned}$$

Comme $D = \text{diag}(d_1, d_2, \dots, d_p)$, on a $D^2 + \lambda I_p = \text{diag}(d_1^2 + \lambda, d_2^2 + \lambda, \dots, d_p^2 + \lambda)$. Cette matrice est clairement inversible comme les d_i sont positifs et $\lambda > 0$. On a donc $(D^2 + \lambda I_p)^{-1} = \text{diag}(\frac{1}{d_1^2 + \lambda}, \frac{1}{d_2^2 + \lambda}, \dots, \frac{1}{d_p^2 + \lambda})$

Cette décomposition est utile quand on souhaite **accélérer le calcul de l'inversion de la matrice** $x^T x + \lambda I_p$.

2.3 Question 3

$$\text{Var}(\beta_{\text{ridge}}^*) = \text{Var}((x^T x + \lambda I_p)^{-1} x^T y) = \text{Var}((x^T x + \lambda I_p)^{-1} x^T (x\beta + \epsilon))$$

Comme supposé au-dessus, x est déterministe donc $\text{Var}(x) = 0$. Par ailleurs, en supposant que $\text{Var}(\epsilon) = \sigma^2 I_n$, on obtient :

$$\begin{aligned}
Var(\beta_{\text{ridge}}^*) &= Var((x^T x + \lambda I_p)^{-1} x^T \epsilon) \\
Var(\beta_{\text{ridge}}^*) &= \sigma^2 (x^T x + \lambda I_p)^{-1} x^T x (x^T x + \lambda I_p)^{-1} \\
Var(\beta_{\text{ridge}}^*) &= \sigma^2 V(D^2 + \lambda I_p)^{-1} V^T V D U^T U D V^T V (D^2 + \lambda I_p)^{-1} V^T \\
Var(\beta_{\text{ridge}}^*) &= \sigma^2 V(D^2 + \lambda I_p)^{-1} D^2 (D^2 + \lambda I_p)^{-1} V^T \\
Var(\beta_{\text{ridge}}^*) &= \sigma^2 \sum_{k=1}^p \frac{d_k^2}{(d_k^2 + \lambda)^2} v_k v_k^T
\end{aligned}$$

$$\text{Or, } Var(\beta_{\text{OLS}}^*) = \sigma^2 V D^{-2} V^T = \sigma^2 \sum_{k=1}^p \frac{1}{d_k^2} v_k v_k^T$$

Comme $\frac{d_k^2}{(d_k^2 + \lambda)^2} < \frac{1}{d_k^2}$, on a $Var(\beta_{\text{ridge}}^*) < Var(\beta_{\text{OLS}}^*)$

2.4 Question 4

Plus λ augmente, plus $Var(\beta_{\text{ridge}}^*)$ diminue. En effet, $Var(\beta_{\text{ridge}}^*) = \sigma^2 \sum_{k=1}^p \frac{d_k^2}{(d_k^2 + \lambda)^2} v_k v_k^T$ est une fonction décroissante de λ .

Regardons à présent le biais :

$$\begin{aligned}
\mathbb{E}[\beta_{\text{Ridge}}^*] &= \mathbb{E}[(x^T x + \lambda I_p)^{-1} x^T y] \\
&= \mathbb{E}[(x^T x + \lambda I_p)^{-1} x^T (x\beta + \epsilon)] \\
&= \mathbb{E}[(x^T x + \lambda I_p)^{-1} (x^T x \beta + x^T \epsilon + \lambda \beta - \lambda \beta)] \\
&= \beta - \lambda (x^T x + \lambda I_p)^{-1} \beta \\
\Rightarrow \text{biais}_\beta(\beta_{\text{Ridge}}^*) &= -\lambda (x^T x + \lambda I_p)^{-1} \beta
\end{aligned}$$

En diagonalisant la matrice $x^T x + \lambda I_p = P(\Delta + \lambda I_p)P^T$, on obtient que $\text{biais}_\beta(\beta_{\text{Ridge}}^*) = -P^T(\frac{\Delta}{\lambda} + I_p)^{-1} P\beta$, donc lorsque λ augmente, le biais tend vers une certaine constante de la forme $-P^T \Gamma_\infty P\beta$.

2.5 Question 5

En supposant que $x^T x = I_p$, on obtient :

$$\begin{aligned}
\beta_{\text{ridge}}^* &= (x^T x + \lambda I_p)^{-1} x^T y \\
&= ((\lambda + 1)I_p)^{-1} x^T y \\
&= \frac{1}{\lambda + 1} x^T y \\
\beta_{\text{OLS}}^* &= (x^T x)^{-1} x^T y = x^T y
\end{aligned}$$

On a donc bien $\beta_{\text{ridge}}^* = \frac{\beta_{\text{OLS}}^*}{\lambda + 1}$ lorsque $x^T x = I_p$.

3 Elastic Net

Par définition, on a :

$$\begin{aligned}
\beta_{ElNet}^* &= \arg \min_{\beta} \|y_c - x_c \beta\|_2^2 + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1 \\
&= \arg \min_{\beta} y_c^T y_c + \beta^T x_c^T x_c \beta - 2\beta^T x_c^T y_c + \lambda_2 \|\beta\|_2^2 + \lambda_1 \sum_{j=1}^d |\beta_j| \\
&= \arg \min_{\beta} \beta^T \beta - 2\beta^T (x_c^T x_c)^{-1} x_c^T y_c + \lambda_2 \|\beta\|_2^2 + \lambda_1 \sum_{j=1}^d |\beta_j| \\
&= \arg \min_{\beta} \|\beta\|_2^2 - 2\beta^T \beta_{OLS}^* + \lambda_2 \|\beta\|_2^2 + \lambda_1 \sum_{j=1}^d |\beta_j| \\
&= \arg \min_{\beta} -2\beta^T \beta_{OLS}^* + (\lambda_2 + 1) \|\beta\|_2^2 + \lambda_1 \sum_{j=1}^d |\beta_j| \\
\beta_{ElNet}^* &= \arg \min_{\beta} \sum_{j=1}^d (-2\beta_j \beta_{OLS,j}^* + (\lambda_2 + 1)\beta_j^2 + \lambda_1 |\beta_j|)
\end{aligned}$$

On peut chercher l'argmin pour chaque terme. On a pour tout j :

$$\begin{aligned}
\beta_{ElNet,j}^* &= \arg \min_{\beta_j} -2\beta_j \beta_{OLS,j}^* + (\lambda_2 + 1)\beta_j^2 + \lambda_1 |\beta_j| \\
&= \arg \min_{\beta_j} \begin{cases} -2\beta_j \beta_{OLS,j}^* + (\lambda_2 + 1)\beta_j^2 + \lambda_1 \beta_j & \text{si } \beta \geq 0 \\ 2\beta_j \beta_{OLS,j}^* + (\lambda_2 + 1)\beta_j^2 - \lambda_1 \beta_j & \text{si } \beta < 0 \end{cases} \\
&= \begin{cases} \frac{2\beta_{OLS,j}^* - \lambda_1}{2(\lambda_2 + 1)} & \text{si } \beta \geq 0 \\ \frac{2\beta_{OLS,j}^* + \lambda_1}{2(\lambda_2 + 1)} & \text{si } \beta < 0 \end{cases} \\
\beta_{ElNet,j}^* &= \frac{\beta_{OLS}^* \pm \frac{\lambda_1}{2}}{1 - \lambda_2}
\end{aligned}$$

4 LDA

Si la covariance de chaque classe est différente, on ne peut plus simplifier par $\frac{1}{|\Sigma|^{1/2}}$ pour toutes les classes. Ainsi, toute la LDA ne peut plus s'appliquer car f^* n'est plus une "**linear** discriminant function". Dans ce cas, il faut utiliser la **QDA**.

$$\begin{aligned}
f^*(x_j) &= \operatorname{argmax}_{C_k} P_{C_k}(x_j) \pi_{C_k} \\
&= \operatorname{argmin}_{C_k} -2\log(P_{C_k}(x_j) \pi_{C_k}) \\
&\propto \operatorname{argmin}_{C_k} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log(|\Sigma_k|) - 2\log(\pi_{C_k}) \\
&= \operatorname{argmin}_{C_k} x^T \Sigma_k^{-1} x - x^T \cdot 2\Sigma_k^{-1} \mu_k + (\mu_k^T \Sigma_k^{-1} \mu_k - 2\log(\pi_{C_k}) + \log(|\Sigma_k|)) \\
&= \operatorname{argmin}_{C_k} x^T c_k x + x^T b_k + a_k
\end{aligned}$$

On voit qu'ici, f^* est une fonction discriminante quadratique (et donc non linéaire). En effet, le terme $x^T c_k x = x^T \Sigma_k^{-1} x$ intervient, i.e. la matrice de covariance de chaque classe intervient.