

Introduction Supervised Learning

IMA205

Pietro Gori

Deadline: Upload the answers and the two notebooks to E-campus. Please verify the deadline on E-Campus.

Theoretical questions

OLS

We have seen that the OLS estimator is equal to $\beta^* = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}$ which can be rewritten as $\beta^* = H\mathbf{y}$. Let $\tilde{\beta} = C\mathbf{y}$ be another linear unbiased estimator of β where C is a $d \times n$ matrix, e.g. $C = H + D$ where D is a non-zero matrix .

- Demonstrate that OLS is the estimator will the smallest variance: compute $\mathbf{E}[\tilde{\beta}]$ and $\text{Var}(\tilde{\beta}) = \mathbf{E}[(\tilde{\beta} - \mathbf{E}[\tilde{\beta}])(\tilde{\beta} - \mathbf{E}[\tilde{\beta}])^T]$ and show when and why $\text{Var}(\beta^*) < \text{Var}(\tilde{\beta})$. Which assumption of OLS do we need to use ?

Ridge regression

Suppose that both \mathbf{y} and the columns of \mathbf{x} are centered (\mathbf{y}_c and \mathbf{x}_c) so that we do not need the intercept β_0 . In this case, the matrix \mathbf{x}_c has d (rather than $d+1$) columns. We can thus write the criterion for ridge regression as:

$$\beta_{ridge}^* = \arg \min_{\beta} (\mathbf{y}_c - \mathbf{x}_c \beta)^T (\mathbf{y}_c - \mathbf{x}_c \beta) + \lambda \|\beta\|_2^2 \quad (1)$$

- Show that the estimator of ridge regression is biased (that is $\mathbf{E}[\beta_{ridge}^*] \neq \beta$).
- Recall that the SVD decomposition is $\mathbf{x}_c = UDV^T$. Write down by hand the solution β_{ridge}^* using the SVD decomposition. When is it useful using this decomposition ? Hint: do you need to invert a matrix ?
- Remember that $\text{Var}(\beta_{OLS}^*) = \sigma^2 (\mathbf{x}^T \mathbf{x})^{-1}$. Show that $\text{Var}(\beta_{OLS}^*) \geq \text{Var}(\beta_{ridge}^*)$.
- When λ increases what happens to the bias and to the variance ? Hint: Compute $\text{MSE} = \mathbf{E}[(y_0 - x_0^T \beta_{ridge}^*)^2]$ at the test point (x_0, y_0) with $y_0 = x_0^T \beta + \epsilon_0$ being the true model and $x_0^T \beta_{ridge}^*$ the ridge estimate.
- Show that $\beta_{ridge}^* = \frac{\beta_{OLS}^*}{1+\lambda}$ when $\mathbf{x}_c^T \mathbf{x}_c = I_d$

OLS

(1)

$$\hat{\beta}^* = Hy \quad \text{where } H = (x^T x)^{-1} x^T$$

OLS estimator

Let $\tilde{\beta} = Cy$ where C is $d \times m$ matrix $C = H + D$, where D
another estimator is non-zero

• Demonstrate that OLS is the estimator with the smallest var

- $E[\beta^*] = H \cdot E[y]$

$$\begin{aligned} \text{Var}(\beta^*) &= E[(\beta^* - E[\beta^*])(\beta^* - E[\beta^*])^T] \\ &= E[H(y - E[y]) (H(y - E[y]))^T] \\ &= H \text{Var}(y) H^T \end{aligned}$$

- $E[\tilde{\beta}] = C \cdot E[y]$

$$\begin{aligned} \text{Var}(\tilde{\beta}) &= E[(\tilde{\beta} - E[\tilde{\beta}])(\tilde{\beta} - E[\tilde{\beta}])^T] \\ &= E[(\tilde{\beta} - C \cdot E[y])(\tilde{\beta} - C \cdot E[y])^T] \\ &= E[(C \cdot y - C \cdot E[y]) \cdot (C \cdot y - C \cdot E[y])^T] \\ &= E[C \cdot (y - E[y]) (C \cdot (y - E[y]))^T] \\ &= E[C \cdot (y - E[y]) (y - E[y])^T \cdot C^T] \\ &= C \cdot E[(y - E[y])(y - E[y])^T] \cdot C^T \end{aligned}$$

$$\text{Var}(\tilde{\beta}) = C \cdot \text{Var}(y) \cdot C^T = (H + D) \text{Var}(y) (H + D)^T$$

$$\text{Var}(\tilde{\beta}) = \underbrace{H \text{Var}(y) H^T}_{\text{Var}(y)} + D \text{Var}(y) D^T$$

$$\text{Var}(\tilde{\beta}) = \text{Var}(\beta^*) + \underbrace{D \text{Var}(y) D^T}_{\text{positive matrix}} > \text{Var}(\beta^*)$$

Ridge Regression

(2)

$$\beta_n^* = \arg \min_{\beta} (y_c - x_c \beta)^T (y_c - x_c \beta) + \lambda \|\beta\|_2^2$$

$$\frac{\partial}{\partial \beta} \left[(y_c - x_c \beta)^T (y_c - x_c \beta) + \lambda \|\beta\|_2^2 \right] = 0$$

$$-2x_c^T (y_c - x_c \beta_n^*) + 2\lambda \beta_n^* = 0$$

$$-2x_c^T y_c + 2x_c^T x_c \beta_n^* + 2\lambda \beta_n^* = 0$$

$$-x_c^T y_c + x_c^T x_c \beta_n^* + \lambda \beta_n^* = 0$$

$\lambda > 0$
always invertible

$$x_c^T y_c = x_c^T x_c \beta_n^* + \lambda \beta_n^* = (x_c^T x_c + \lambda I) \beta_n^*$$

$$\beta_n^* = (x_c^T x_c + \lambda I)^{-1} x_c^T y_c //$$

$$\mathbb{E}[\beta_n^*] = \mathbb{E}\left[-(x_c^T x_c + \lambda)^{-1} x_c^T y_c\right] = (x_c^T x_c + \lambda)^{-1} x_c^T \mathbb{E}[y_c] //$$

$$\mathbb{E}[\beta_n^*] \neq \beta$$

$$③ \quad x_c = UDV^T$$

↳ From the previous question

$$\beta_n^* = (x_c^T x_c + \lambda I)^{-1} x_c^T y_c$$

$$\beta_n^* = [(UDV^T)^T (UDV^T) + \lambda I]^{-1} (UDV^T)^T y_c$$

$$\beta_n^* = [V D^T U^T U D V^T + \lambda I]^{-1} V D^T U^T y_c$$

Since D is a diagonal matrix $D^T = D$

$$\beta_n^* = [V D^2 V^T + \lambda I]^{-1} V D U^T y_c$$

Since V is an orthonormal matrix $V^T = V^{-1} \Rightarrow \underline{VV^T = I}$

$$\beta_n^* = [V D^2 V^T + \lambda I V V^T]^{-1} V D U^T y_c$$

$$\beta_n^* = [V (D^2 + \lambda I) V^T]^{-1} V D U^T y_c$$

$$\beta_n^* = [V^T (D^2 + \lambda I)]^{-1} V^{-1} V D U^T y_c$$

$$\boxed{\beta_n^* = V (D^2 + \lambda I)^{-1} D U^T y_c}$$

↳ Since D is a diagonal matrix, $D^2 + \lambda I$ is also a diagonal matrix. It means that the inverse matrix is a diagonal matrix with elements $e_{d,i} = \frac{1}{e_d^2 + \lambda}$, where e_d are the elements from the D matrix. So, it's not more necessary calculate the inverse matrix which could be a computational costly operation.

$$\textcircled{4} \quad \text{Var}(\beta_{OLS}^*) = \sigma^2 (X^T X)^{-1},$$

Show that $\text{Var}(\beta_n^*) \leq \text{Var}(\beta_{OLS}^*)$

$$\beta_n^* = (X_c^T X_c + \lambda I)^{-1} X_c^T y_c$$

$$E[\beta_n^*] = (X_c^T X_c + \lambda I)^{-1} X_c^T E[y_c]$$

$$\text{Var}(\beta_n^*) = E[(\beta_n^* - E[\beta_n^*])(\beta_n^* - E[\beta_n^*])^T]$$

→ By means of notation let's assume $H = (X_c^T X_c + \lambda I)^{-1} X_c^T$

So,

$$\text{Var}(\beta_n^*) = E[(H y_c - H E[y_c])(H y_c - H E[y_c])^T]$$

$$= E[H(y_c - E[y_c])(y_c - E[y_c])^T H^T]$$

$$= H \cdot \underbrace{E[(y_c - E[y_c])(y_c - E[y_c])^T]}_{\text{Var}(y) = \sigma^2} H^T$$

$$\text{Var}(y) = \sigma^2$$

$$\Leftrightarrow \text{Var}(\beta_n^*) = \sigma^2 H H^T = \sigma^2 [(X_c^T X_c + \lambda I)^{-1} X_c^T][(X_c^T X_c + \lambda I)^{-1} X_c^T]^T$$

$$\text{Var}(\beta_n^*) = \sigma^2 (X_c^T X_c + \lambda I)^{-1} X_c^T X_c [(X_c^T X_c + \lambda I)^{-1}]^T$$

Since $\lambda > 0$ and $X_c^T X_c > 0$, we can say that

$$(X_c^T X_c)^{-1} \left[(X_c^T X_c)^{-1} \right]^T \geq (X_c^T X_c + \lambda I)^{-1} \left[(X_c^T X_c + \lambda I)^{-1} \right]^T$$

So, is also true:

$$(x_c^T x_c)^{-1} x_c^T x_c [(x_c^T x_c)^{-1}]^T \geq (x_c^T x_c + \lambda I)^{-1} x_c^T x_c [(x_c^T x_c + \lambda I)^{-1}]$$

$$[(x_c^T x_c)^{-1}]^T \geq (x_c^T x_c + \lambda I)^{-1} x_c^T x_c [(x_c^T x_c + \lambda I)^{-1}]$$

Also true for the transpose matrix

$$(x_c^T x_c)^{-1} \geq (x_c^T x_c + \lambda I)^{-1} x_c^T x_c [(x_c^T x_c + \lambda I)^{-1}]^T$$

$$\sigma^2 (x_c^T x_c)^{-1} \geq \sigma^2 (x_c^T x_c + \lambda I)^{-1} x_c^T x_c [(x_c^T x_c + \lambda I)^{-1}]^T$$

Finally

$$\text{Var}(\beta_{OLS}^*) \geq \text{Var}(\beta_R^*) //$$

(5) $\lambda \rightarrow \infty$ $\text{Var}(\beta_n^*) \rightarrow ?$ $\text{Bias}(\beta_n^*) \rightarrow ?$

$$\text{Var}(\beta_n^*) = \underline{\sigma^2} \underline{(X_C^T X_C + \lambda I)^{-1}} X_C^T X_C \underline{[(X_C^T X_C + \lambda I)^{-1}]}^T$$

\hookrightarrow independent from lambda

When $\lambda \rightarrow \infty$ $\text{Var}(\beta_n^*) \rightarrow 0$

For $\lambda \rightarrow \infty$ $\beta_n^* \rightarrow 0$,

$$MSE = E[(y_0 - X_0^T \beta_n^*)^2]$$

$$y_0 = X_0^T \beta + \varepsilon$$

$$MSE = E[(X_0^T (\beta - \beta_n^*) + \varepsilon_0)^2]$$

$$MSE = E[(X_0^T (\beta - \beta_n^*))^2 + 2(X_0^T (\beta - \beta_n^*)) \varepsilon_0 + \varepsilon_0^2]$$

\hookrightarrow Let's assume $\varepsilon_0 \sim N(0, \text{Var}(\varepsilon_0))$

$$MSE = (X_0^T (\beta - \beta_n))^2 + 2(X_0^T (\beta - \beta_n^*)) E[\varepsilon_0] + \text{Var}(\varepsilon_0)$$

$$MSE = \underbrace{(X_0^T (\beta - \beta_n))^2}_{\text{bias}^2} + \text{Var}(\varepsilon_0)$$

With $\lambda \rightarrow \infty$, $\text{Var} \rightarrow 0$, $\beta_n \rightarrow 0$ which means that the bias term increases.

⑥ Show that $\beta_n^* = \frac{\beta_{OLS}}{1 + \lambda}$ when $X_c^T X_c = I$

$$\beta_n^* = (X_c^T X_c + \lambda I)^{-1} X_c^T y_c \quad \beta_{OLS} = (X_c^T X_c)^{-1} X_c^T y_c$$

$$X_c^T X_c = I \Rightarrow \beta_{OLS} = X_c^T y_c$$

$$\Rightarrow \beta_n^* = (I + \lambda I)^{-1} X_c^T y_c$$

$$\boxed{\beta_n^* = \frac{\beta_{OLS}}{1 + \lambda}} //$$

Elastic Net

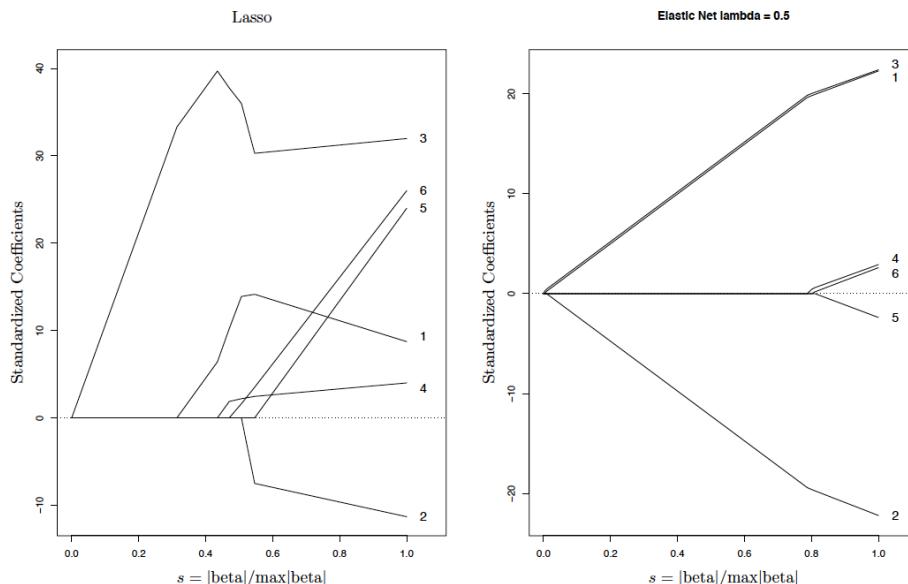
Using the previous notation, we can also combine Ridge and Lasso in the so-called Elastic Net regularization :

$$\beta_{ElNet}^* = \arg \min_{\beta} (\mathbf{y}_c - \mathbf{x}_c \beta)^T (\mathbf{y}_c - \mathbf{x}_c \beta) + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1 \quad (2)$$

Calling $\alpha = \frac{\lambda_2}{\lambda_1 + \lambda_2}$, solving the previous Eq. is equivalent to:

$$\beta_{ElNet}^* = \arg \min_{\beta} (\mathbf{y}_c - \mathbf{x}_c \beta)^T (\mathbf{y}_c - \mathbf{x}_c \beta) + \lambda \left[\alpha \left(\sum_{j=1}^d \beta_j^2 \right) + (1 - \alpha) \left(\sum_{j=1}^d |\beta_j| \right) \right] \quad (3)$$

- This regularization overcomes some of the limitations of the Lasso, notably:
 - If $d > N$ Lasso can select at most N variables → ElNet removes this limitation
 - If a group of variables are highly correlated, Lasso randomly selects only one variable → with ElNet correlated variables have a similar value (grouped)
 - Lasso solution paths tend to vary quite drastically → ElNet regularizes the paths
 - If $N > d$ and there is high correlation between the variables, Ridge tends to have a better performance in prediction → ElNet combines Ridge and Lasso to have better (or similar) prediction accuracy with less (or more grouped) variables



- Compute by hand the solution of Eq.2 supposing that $\mathbf{x}_c^T \mathbf{x}_c = I_d$ and show that the solution is: $\beta_{ElNet}^* = \frac{(\beta_{OLS}^*)_j \pm \frac{\lambda_1}{2}}{1 + \lambda_2}$

Elastic Net

$$\textcircled{7} \quad X_C^T X_C = I \quad \beta_E^* = \frac{\beta_{OLS}^* \pm \frac{\lambda_2}{2}}{1 + \frac{\lambda_2}{2}} \quad \beta_{OLS} = (X_C^T X_C)^{-1} X_C^T y_C \\ = X_C^T y_C$$

$$\beta_i^* = \underset{\beta}{\operatorname{arg\,min}} (y_C - X_C \beta)^T (y_C - X_C \beta) + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1$$

$$\frac{\partial}{\partial \beta_i^*} = -2X_C^T (y_C - X_C \beta_E^*) + 2\lambda_2 \beta_E^* + \lambda_1 \underbrace{\begin{cases} -1, & \beta < 0 \\ [-1, 1], & \beta = 0 \\ 1, & \beta > 1 \end{cases}}_{\text{subgradient}} = 0$$

21.1

$$\beta_{OLS}^* \quad I \\ -2X_C^T y_C - 2X_C^T X_C \beta_E^* + 2\lambda_2 \beta_E^* \pm \lambda_1 = 0$$

$$-2\beta_{OLS}^* - 2\beta_E^* (1 - \lambda_2) \pm \lambda_1 = 0$$

$$\beta_{OLS}^* + \beta_E^* (1 - \lambda_2) \mp \frac{\lambda_1}{2} = 0$$

$$\beta_E^* (1 - \lambda_2) = -\beta_{OLS}^* \mp \frac{\lambda_1}{2} = 0$$

$$\beta_E^* = \frac{\beta_{OLS}^* \pm \frac{\lambda_1}{2}}{\lambda_2 - 1} //$$