

IMA205 TP3 Theoretical Questions

Felipe Vicentin

March 14, 2025

OLS

Questions

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

Answers

1. Let $\tilde{\beta} = (H + D)\mathbf{y}$. Then, using the fact that $\text{Var}(A\mathbf{x}) = A\text{Var}(\mathbf{x})A^\top$ for unbiased \mathbf{x} and that $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$, we have:

$$\begin{aligned}\text{Var}(\tilde{\beta}) &= \text{Var}((H + D)\mathbf{y}) \\ &= \text{Var}(H\mathbf{y}) + \text{Var}(D\mathbf{y}) \\ &= \text{Var}(\beta^*) + D\text{Var}(\mathbf{y})D^\top \\ &= \text{Var}(\beta^*) + D(\mathbf{y}\mathbf{y}^\top)D^\top \\ &= \text{Var}(\beta^*) + \|D\mathbf{y}\|_2^2\end{aligned}$$

It is evident that $\text{Var}(\tilde{\beta}) \geq \text{Var}(\beta^*)$. Strict inequality is obtained when the spectral radius of D is non-zero (i.e., $D \neq 0$). Of course, we are assuming the OLS estimator is unbiased. This is easily verifiable when considering that $(X^\top X)^{-1}$ exists and that the error in the data has 0 mean:

$$\begin{aligned}\mathbf{E}[\beta^*] &= \mathbf{E}[(X^\top X)^{-1}X^\top \mathbf{y}] \\ &= \mathbf{E}[(X^\top X)^{-1}X^\top (X\beta + \varepsilon)] \\ &= (X^\top X)^{-1}X^\top (\mathbf{E}[X\beta] + \mathbf{E}[\varepsilon]) \\ &= \beta\end{aligned}$$

Ridge regression

Questions

2. Show that the estimator of ridge regression is biased (that is $\mathbf{E}[\beta_{\text{ridge}}^*] \neq \beta$).
3. Recall that the SVD decomposition is $\mathbf{x}_c = UDV^\top$. 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?
4. Remember that $\text{Var}(\beta_{\text{OLS}}^*) = \sigma^2(\mathbf{x}^\top \mathbf{x})^{-1}$. Show that $\text{Var}(\beta_{\text{OLS}}^*) \geq \text{Var}(\beta_{\text{ridge}}^*)$.
5. When λ increases what happens to the bias and to the variance? Hint: Compute $\text{MSE} = \mathbf{E}[(y_0 - x_0^\top \beta_{\text{ridge}}^*)^2]$ at the test point (x_0, y_0) with $y_0 = x_0^\top \beta + \epsilon_0$ being the true model and $x_0^\top \beta_{\text{ridge}}^*$ the ridge estimate.
6. Show that $\beta_{\text{ridge}}^* = \frac{\beta_{\text{OLS}}^*}{1+\lambda}$ when $\mathbf{x}_c^\top \mathbf{x}_c = I_d$.

Answers

2. As seen in class, $\beta_{\text{ridge}}^* = (X^\top X + \lambda I)^{-1} X^\top \mathbf{y}$. Then, it is easy to see that

$$\begin{aligned}\mathbf{E}[\beta_{\text{ridge}}^*] &= \mathbf{E}[(X^\top X + \lambda I)^{-1} X^\top \mathbf{y}] \\ &= (X^\top X + \lambda I)^{-1} X^\top (\mathbf{E}[X\beta] + \mathbf{E}[\varepsilon]) \\ &= (X^\top X + \lambda I)^{-1} X^\top X\beta\end{aligned}$$

Since $(X^\top X + \lambda I)^{-1} X^\top X \neq I$ for $\lambda > 0$, the estimator is biased.

3. If we use $X = UDV^\top$, we have

$$\begin{aligned}\beta_{\text{ridge}}^* &= (X^\top X + \lambda I)^{-1} X^\top \mathbf{y} \\ &= ((UDV^\top)^\top UDV^\top + \lambda I)^{-1} (UDV^\top)^\top \mathbf{y} \\ &= (VDU^\top UDV^\top + \lambda I)^{-1} VDU^\top \mathbf{y} \\ &= (V(D^2 + \lambda I)V^\top)^{-1} VDU^\top \mathbf{y} \\ &= V(D^2 + \lambda I)^{-1} V^\top VDU^\top \mathbf{y} \\ &= V(D^2 + \lambda I)^{-1} DU^\top \mathbf{y}\end{aligned}$$

We see that $D^2 + \lambda I$ is diagonal and no inversion is needed. This speeds up the algorithm substantially. The form above is particularly useful when testing many λ to find the best hyperparameter.

4. Let us first compute the variance of the ridge estimator.

$$\begin{aligned}\text{Var}(\beta_{\text{ridge}}^*) &= \text{Var}((X^\top X + \lambda I)^{-1} X^\top \mathbf{y}) \\ &= (X^\top X + \lambda I)^{-1} X^\top \text{Var}(\mathbf{y}) X (X^\top X + \lambda I)^{-1} \\ &= (X^\top X + \lambda I)^{-1} X^\top \text{Var}(\varepsilon) X (X^\top X + \lambda I)^{-1} \\ &= \sigma^2 (X^\top X + \lambda I)^{-1} X^\top X (X^\top X + \lambda I)^{-1}\end{aligned}$$

We know that the eigenvalues of $X^\top X + \lambda I$ are equal to the eigenvalues of $X^\top X$ plus $\lambda \geq 0$, so $X^\top X \preceq X^\top X + \lambda I$. In turn, this implies that:

$$\begin{aligned}(X^\top X + \lambda I)^{-1} &\preceq (X^\top X)^{-1} \\ (X^\top X + \lambda I)^{-1} (X^\top X) &\preceq (X^\top X)^{-1} (X^\top X) \\ (X^\top X + \lambda I)^{-1} (X^\top X) (X^\top X + \lambda I)^{-1} &\preceq (X^\top X)^{-1} (X^\top X) (X^\top X + \lambda I)^{-1} \\ (X^\top X + \lambda I)^{-1} (X^\top X) (X^\top X + \lambda I)^{-1} &\preceq (X^\top X)^{-1} (X^\top X) (X^\top X)^{-1} \\ \sigma^2 (X^\top X + \lambda I)^{-1} (X^\top X) (X^\top X + \lambda I)^{-1} &\preceq \sigma^2 (X^\top X)^{-1} \\ \text{Var}(\beta_{\text{ridge}}^*) &\preceq \text{Var}(\beta_{\text{OLS}}^*)\end{aligned}$$

5. The bias of the Ridge estimator is given by

$$\begin{aligned}b(\beta_{\text{ridge}}^*) &= \mathbf{E}(\beta_{\text{ridge}}^*) - \beta \\ &= ((X^\top X + \lambda I)^{-1} X^\top X - I)\beta\end{aligned}$$

We see that as λ increases, the bias gets greater in absolute value. The variance, on the other hand, gets smaller.

6. If $X^\top X = I$, then $\beta_{\text{OLS}}^* = (X^\top X)^{-1} X^\top \mathbf{y} = X^\top \mathbf{y}$. The ridge estimator is $\beta_{\text{ridge}}^* = (X^\top X + \lambda I)^{-1} X^\top \mathbf{y} = (1 + \lambda)^{-1} X^\top \mathbf{y} = \frac{1}{1 + \lambda} X^\top \mathbf{y}$. So, we can conclude that $\beta_{\text{ridge}}^* = \frac{\beta_{\text{OLS}}^*}{1 + \lambda}$.

Elastic Net

Questions

7. 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:

$$(\beta_{\text{ElNet}}^*)_j = \frac{(\beta_{\text{OLS}}^*)_j \pm \frac{\lambda_1}{2}}{1 + \lambda_2}$$

Answers

7.

$$\begin{aligned} f(\beta) &= (\mathbf{y}_c - \mathbf{x}_c \beta)^T (\mathbf{y}_c - \mathbf{x}_c \beta) + \lambda_2 \|\beta\|_2^2 \lambda_1 \|\beta\|_1 \\ &= \mathbf{y}_c^T \mathbf{y}_c - 2\mathbf{y}_c^T \mathbf{x}_c \beta + \beta^T \mathbf{x}_c^T \mathbf{x}_c \beta + \lambda_2 \beta^T \beta + \lambda_1 \|\beta\|_1 \\ &= \mathbf{y}_c^T \mathbf{y}_c - 2\mathbf{y}_c^T \mathbf{x}_c \beta + (1 + \lambda_2) \beta^T \beta + \lambda_1 \|\beta\|_1 \\ \implies \partial f(\beta) &= -2\mathbf{x}_c^T \mathbf{y}_c + 2(1 + \lambda_2)\beta + \lambda_1 \partial(\|\cdot\|_1)(\beta) \\ \implies -\mathbf{x}_c^T \mathbf{y}_c + (1 + \lambda_2)\beta^* &= -\frac{\lambda_1}{2} \partial(\|\cdot\|_1)(\beta^*) \\ \implies \beta_j^* &\in \begin{cases} \left\{ \frac{\mathbf{x}_c^T \mathbf{y}_c + \frac{\lambda_1}{2}}{\lambda_2 + 1} \right\} & \text{if } \beta_j^* < 0 \\ \left\{ \frac{\mathbf{x}_c^T \mathbf{y}_c - \frac{\lambda_1}{2}}{\lambda_2 + 1} \right\} & \text{if } \beta_j^* > 0 \\ \left[\frac{\mathbf{x}_c^T \mathbf{y}_c - \frac{\lambda_1}{2}}{\lambda_2 + 1}, \frac{\mathbf{x}_c^T \mathbf{y}_c + \frac{\lambda_1}{2}}{\lambda_2 + 1} \right] & \text{if } \beta_j^* = 0 \end{cases} \end{aligned}$$

Since $\beta_{\text{OLS}}^* = \mathbf{x}_c^T \mathbf{y}_c$, we have that:

$$(\beta_{\text{ElNet}}^*)_j = \begin{cases} \frac{(\beta_{\text{OLS}}^*)_j - \frac{\lambda_1}{2}}{\lambda_2 + 1} & (\beta_{\text{OLS}}^*)_j > \frac{\lambda_1}{2} \\ \frac{(\beta_{\text{OLS}}^*)_j + \frac{\lambda_1}{2}}{\lambda_2 + 1} & (\beta_{\text{OLS}}^*)_j < -\frac{\lambda_1}{2} \\ 0 & |(\beta_{\text{OLS}}^*)_j| \leq \frac{\lambda_1}{2} \end{cases}.$$