

SDS 387 Linear Models

Fall 2024

Lecture 18 - Thu, Oct 31, 2024

Instructor: Prof. Ale Rinaldo

GAUSS MARKOV THEOREM

We are in the setting where Φ is a deterministic design matrix and the linear model is well-specified: $\text{rank}(\Phi) = d < n$

$$Y = \Phi \beta^* + \varepsilon$$

$\varepsilon \sim (0, \sigma^2 I_n)$

Any estimator of the form $A Y$, where A is a deterministic matrix (possibly depending on Φ) is said to be linear and is unbiased if

$$E[A Y] = \beta^*, \quad \forall \beta^*$$

We know that $\hat{\beta} = (\Phi^T \Phi)^{-1} \Phi^T Y$ is an unbiased linear estimator of β^* .

GAUSS-MARKOV THEOREM: $\hat{\beta}$ is the best linear unbiased estimator (BLUE) of β^*

This means that, among all unbiased linear estimators of β^* of the form $\underset{d \times d}{A}Y$,

$$\text{Var}[\hat{\beta}] \leq \text{Var}[AY]$$

where we say that a $d \times d$ matrix B is ≤ 0 (non-negative in the positive-semidefinite or Loewner partial order) if B is posd. Then we can write $\underset{d \times d}{A} \leq \underset{d \times d}{B}$ to mean that $A - B \geq 0$.

Remark This is a partial order, meaning that it is possible that A and B are not comparable meaning neither $A \leq B$ nor $B \leq A$ holds. In particular $A \not\leq B$ does not mean $A \geq B$.

PF/ Let AY be an unbiased estimator of β^* . Then

$$\beta^* = \mathbb{E}[AY] = A\Phi\beta^* + \underbrace{A\mathbb{E}[\varepsilon]}_{=0} = A\Phi\beta^*$$

\hookrightarrow this is true for any β^* so $A\Phi = I_{d \times d}$

Now $\text{Var}[AY] = \sigma^2 AA^T$ (because $\text{Var}[Y] = \text{Var}[\varepsilon] = \sigma^2 I_n$)

Next let $D = A - (\Phi^T \Phi)^{-1} \Phi^T$. So

$$\text{Var}[AY] = \sigma^2 (D + (\Phi^T \Phi)^{-1} \Phi^T) (D + (\Phi^T \Phi)^{-1} \Phi^T)^T$$

Because $A\Phi = I$, $(D + (\Phi^T \Phi)^{-1} \Phi^T)\Phi = I$ or

$$D\Phi + I = I \Rightarrow D\Phi = 0$$

As a result

$$\text{Var}[AY] = \sigma^2 DD^T + \sigma^2 (\Phi^T \Phi)^{-1} \Phi^T \Phi (\Phi^T \Phi)^{-1} \\ + \text{cross terms that are zero}$$

$$= \underbrace{\sigma^2 DD^T}_{\geq 0} + \text{Var}[\hat{\beta}]$$

↓

$$\text{Var}[AY] - \text{Var}[\hat{\beta}] \geq 0 \quad \text{or} \quad \text{Var}[AY] \geq \text{Var}[\hat{\beta}]$$

In fact $c^T \hat{\beta}$ for any $c \in \mathbb{R}^d$ is the BLUE of $c^T \beta^*$

RIDGE REGRESSION

Suppose that d is large compared to n , meaning that $\frac{d}{n}$ is close to 1. There are of course computational issues because $\Phi^T \Phi$ is not well conditioned and statistical issues because $\text{Var}[\hat{\beta}] = \sigma^2 (\Phi^T \Phi)^{-1} \rightarrow$ poorly conditioned

One approach is to regularize, to solve a penalized least-squares problem that encourages "good" properties of the solutions. Ridge regression is an important example (and so is Lasso).

Note: here we are interested in minimizing the prediction risk or the excess risk.

The ridge least squares estimator is

$$\hat{\beta}_\lambda = \underset{\beta \in \mathbb{R}^d}{\text{argmin}} \frac{\|y - \Phi \beta\|^2}{n} + \lambda \|\beta\|^2 \quad \lambda > 0$$

↓
regularization parameter

if $\lambda = 0$ this gives us $\hat{\beta}_0$, the OLS

$\hat{\beta}_\lambda$ is uniquely defined no matter Φ and is equal to

$$\hat{\beta}_\lambda = \underbrace{\left(\frac{\Phi^T \Phi}{n} + \lambda I_d \right)^{-1}}_{\text{always invertible}} \frac{\Phi^T y}{n} = \left(\frac{\lambda}{n} + \lambda I_d \right)^{-1} \frac{\Phi^T y}{n}$$

PP/ Let $F_\lambda(\beta) = \frac{1}{n} \|\mathcal{Y} - \Phi\beta\|^2 + \lambda \|\beta\|^2$ is strictly convex so $\hat{\beta}_\lambda$ is found by checking first order optimality conditions:

$$0 = \nabla F_\lambda(\hat{\beta}_\lambda) = \frac{2}{n} \Phi^T (\Phi \hat{\beta}_\lambda - \mathcal{Y}) + 2\lambda \hat{\beta}_\lambda \quad \square$$

Remarks *hw!*

1) $\lim_{\lambda \rightarrow 0} \hat{\beta}_\lambda = \hat{\beta}_{MN}$ *hw*

2) Alternative expression:

$$\hat{\beta}_\lambda = \frac{\Phi^T}{n} \left(\underbrace{\Phi\Phi^T}_{n \times n} + \lambda I_n \right)^{-1} \mathcal{Y}$$

better numerically if $\text{rank}(\Phi) = n < d$

3) Let $\Phi = U \Sigma V^T$. Then

$$\Phi \hat{\beta}_\lambda = \sum_{i=1}^{\text{rank}(\Phi)} \mu_i \langle \mathcal{Y}, \mu_i \rangle \frac{\sigma_i^2}{\sigma_i^2 + \lambda}$$

μ_i \leftarrow i th column of U

compare to

$$\Phi \hat{\beta}_{MN} = \sum_{i=1}^{\text{rank}(\Phi)} \mu_i \langle \mathcal{Y}, \mu_i \rangle$$

σ_i \leftarrow i th singular value

Statistical properties of $\hat{\beta}_\lambda$ (recall we want to minimize the risk!)

Proposition 3.7 The excess risk of $\hat{\beta}_\lambda$

$$\mathbb{E}[R(\hat{\beta}_\lambda)] - R(\beta^*) = \lambda^2 \beta^{*\top} \left(\sum_{i=1}^{\text{rank}(\Phi)} \frac{1}{\sigma_i^2} + \lambda I_d \right)^{-2} \sum_{i=1}^{\text{rank}(\Phi)} \beta_i^{*2} +$$

$$\frac{\sigma^2}{n} \text{tr} \left(\hat{\Sigma}^{-1} (\hat{\Sigma} + \lambda \mathbf{I}_d)^{-1} \right)$$

$$= B + V$$

B is a bias term increasing in λ and V is the variance term, decreasing in λ

Note if $\hat{\Sigma}$ is invertible, when $\lambda \rightarrow 0$ we receive the risk of $\hat{\beta}$ (which is $\sigma^2 \frac{d}{n}$).

PA/ Recall the decomposition of the excess risk of any estimator $\tilde{\beta}$:

$$\mathbb{E}[R(\tilde{\beta})] - R(\beta^*) = \underbrace{\|\mathbb{E}[\tilde{\beta}] - \beta^*\|_{\hat{\Sigma}}^2}_{\sigma^2} + \mathbb{E}[\|\tilde{\beta} - \mathbb{E}[\tilde{\beta}]\|_{\hat{\Sigma}}^2]$$

Now replace $\tilde{\beta}$ by $\hat{\beta}_\lambda$. For the bias:

$$\begin{aligned} \mathbb{E}[\hat{\beta}_\lambda] &= (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \frac{\Phi^T}{n} \Phi \beta^* + (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \frac{\Phi^T}{n} \mathbb{E}[\varepsilon] \\ &= (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \hat{\Sigma} \beta^* \\ &= \beta^* - \lambda (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \beta^* \end{aligned}$$

Exercise

$$\begin{aligned} \hookrightarrow \|\mathbb{E}[\hat{\beta}_\lambda] - \beta^*\|_{\hat{\Sigma}}^2 &= \lambda^2 \beta^{*\top} (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda \mathbf{I})^{-1} \beta^* \\ &= \lambda^2 \beta^{*\top} (\hat{\Sigma} + \lambda \mathbf{I}_d)^{-2} \hat{\Sigma} \beta^* \end{aligned}$$

because $(\hat{\Sigma} + \lambda \mathbf{I})^{-1}$ and $\hat{\Sigma}$ commute. HW

As for the variance term:

$$\begin{aligned}
 \mathbb{E} \left[\|\hat{\beta}_b - \mathbb{E}[\hat{\beta}_b]\|^2_{\hat{\Sigma}} \right] &= \mathbb{E} \left[\left\| (\hat{\Sigma} + \lambda I_d)^{-1} \frac{\Phi^T \epsilon}{n} \right\|_{\hat{\Sigma}}^2 \right] \\
 &= \frac{1}{n^2} \mathbb{E} \left[\epsilon^T \Phi (\hat{\Sigma} + \lambda I)^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda I)^{-1} \Phi^T \epsilon \right] \\
 &= \frac{1}{n^2} \mathbb{E} \left[\text{tr} \left(\Phi^T \epsilon \epsilon^T \Phi (\hat{\Sigma} + \lambda)^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda I)^{-1} \right) \right] \\
 &= \frac{\sigma^2}{n} \text{tr} \left(\hat{\Sigma} (\hat{\Sigma} + \lambda)^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda)^{-1} \right) \\
 &= \frac{\sigma^2}{n} \text{tr} \left(\hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2} \right) \\
 &= \sum_{i=1}^p \frac{\hat{d}_i^2}{(\hat{d}_i + \lambda)^2} \quad \begin{array}{l} \hat{d}_i \text{ } i\text{th eigenvalue} \\ \text{of } \hat{\Sigma} \end{array}
 \end{aligned}$$

degrees of freedom