

Lecture 20, Thu Nov 6

• GAUSS MARKOV THEOREM

Assume a linear model with fixed covariates, i.e. $\Phi \in \mathbb{R}^{n \times d}$ is deterministic:

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

\hookrightarrow iid errors $\sim (0, \sigma^2)$

\downarrow

deterministic homoschedastic errors

Any estimation of β^* of the form $A Y$ is a linear unbiased estimator of β^* if

$$\mathbb{E}[AY] = \beta^* \quad (\text{for all } \beta^* \in \mathbb{R}^d)$$

$\hookrightarrow \hat{\beta} = (\Phi^\top \Phi)^{-1} \Phi^\top Y$ is a linear unbiased estimator.

• How good is $\hat{\beta}$ compared to all linear unbiased estimators?

Gauss-Markov Thm $\hat{\beta}$ is the $\xrightarrow{\text{Best Linear Unbiased Estimator}} \text{BLUE}!$

Remark: "Best" means that, for any other linear unbiased estimator AY , $\text{Var}[\hat{\beta}] \leq \text{Var}[AY]$

where for 2 pos matrices M_1 and M_2 of the same size

$$M_1 \preceq M_2 \iff M_2 - M_1 \succeq 0$$

\nwarrow pos. order $\hookrightarrow = \text{is positive semidefinite}$

This means that, for any $x \in \mathbb{R}^d$,

$$x^\top \text{Var}[\tilde{\beta}] x \leq x^\top \text{Var}[AY] x$$

choose x to be the i^{th} standard basis vector $e_i = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} \rightarrow i^{\text{th}} \text{ position}$

to conclude that

the i^{th} element of diagonal of $\text{Var}[\tilde{\beta}]$ is

smaller than i^{th} element $\text{Var}[AY]$

$$\hookrightarrow \text{Var}[\tilde{\beta}_i] \leq \text{Var}[\tilde{\beta}_i] \quad \tilde{\beta} = AY$$

In fact this implies

$$\text{Var}[c^\top \tilde{\beta}] \leq \text{Var}[c^\top \beta^*] \quad \text{for all } c \in \mathbb{R}^d$$

\sim
contrast

Pf/ Let $\tilde{\beta} = AY$ s.t. $E[\tilde{\beta}] = \beta^*$. Then

$$\beta^* = E[AY] = E[A\Phi\beta^* + A\epsilon] = A\Phi\beta^* + A\epsilon$$

\Downarrow

$$A\Phi = I_{\text{at risk}}$$

$$\text{Next } \text{Var}[\tilde{\beta}] = A \underbrace{\text{Var}[Y]}_{6^2 I_n} A^\top = 6^2 A A^\top$$

$$\text{Let } D = A - (\Phi^\top \Phi)^{-1} \Phi^\top \quad \text{so}$$

$$\text{Var}[\tilde{\beta}] = \text{Var}[AY] = 6^2 (D + (\Phi^\top \Phi)^{-1} \Phi^\top) (D + (\Phi^\top \Phi)^{-1} \Phi^\top)^\top$$

$$\text{But } \underbrace{(D + (\Phi^\top \Phi)^{-1} \Phi^\top)}_A \Phi = I_d \quad \text{so}$$

$$\text{Var}[\hat{\beta}] = \underbrace{\sigma^2 D D^\top}_{\geq 0} + \underbrace{\sigma^2 (\Phi^\top \Phi)^{-1} \Phi^\top \Phi (\Phi^\top \Phi)^{-1}}_{\sigma^2 (\Phi^\top \Phi)^{-1}} = \text{Var}[\hat{\beta}]$$

↓

$$\text{Var}[\tilde{\beta}] \geq \text{Var}[\hat{\beta}] \quad \blacksquare$$

- ffw : extension to random Φ !

RIDGE REGRESSION

- Chapter 3 for Bach's book Learning theory from first principles
- Suppose d is large compared to n , almost close to n
 in this case $\Phi^\top \Phi$ may not be well-conditioned
 i.e. $\frac{\lambda_{\max}(\Phi^\top \Phi)}{\lambda_{\min}(\Phi^\top \Phi)}$ is large and $(\Phi^\top \Phi)^{-1}$ is unstable
- $\text{Var}[\hat{\beta}] = \sigma^2 (\Phi^\top \Phi)^{-1}$ is unstable
- Our approach is to regularize - to solve a penalized least squares problem that includes a penalty for not choosing "good" solutions.
- The ridge regression estimator arises as the solution to this problem:

For a value $\lambda \geq 0$, let
 \hookrightarrow penalty parameter

ridge estimator $\leftarrow \hat{\beta}_\lambda = \underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} \underbrace{\frac{\|\mathbf{y} - \underline{\Phi}\beta\|^2}{n} + \lambda \|\beta\|^2}_{f(\beta)}$

$$= \left(\frac{\underline{\Phi}^\top \underline{\Phi}}{n} + \lambda I_d \right)^{-1} \frac{\underline{\Phi}^\top \mathbf{y}}{n} = \left(\Sigma + \lambda I_d \right)^{-1} \frac{\underline{\Phi}^\top \mathbf{y}}{n}$$

when $\lambda = 0$ this reduces to OLS $\hat{\beta}$ (if Σ is invertible)

or the min-norm estimator
 $\underline{\Phi}^\top \mathbf{y}$ if it is not

Pf/ The objective function is strictly convex so the first order optimality conditions are

$$0 = \nabla f(\hat{\beta}_\lambda) = \frac{2}{n} \underline{\Phi}^\top (\underline{\Phi}\hat{\beta}_\lambda - \mathbf{y}) + 2\lambda \hat{\beta}_\lambda$$

\hookrightarrow solve for $\hat{\beta}_\lambda$ \equiv

This solution exists, is unique even if $d > n$ and regardless of how poorly conditioned $\underline{\Phi}\underline{\Phi}^\top$ is!

Remarks

1) if $\lambda \rightarrow 0$ then

$$\hat{\beta}_\lambda \rightarrow \hat{\beta}_{\text{min-norm}} \text{ or } \hat{\beta} \text{ if } \underline{\Phi}^\top \underline{\Phi} \text{ is invertible}$$

2) Alternative expression:

$$\hat{\beta}_\lambda = \left(\frac{\underline{\Phi}^\top \underline{\Phi}}{n} + \lambda I_d \right)^{-1} \frac{\underline{\Phi}^\top \mathbf{y}}{n} = \frac{\underline{\Phi}^\top}{n} \left(\frac{\underline{\Phi} \underline{\Phi}^\top}{n} + \lambda I_d \right)^{-1} \mathbf{y}$$

3) Let $\underline{\Phi} = U \Sigma V^T$ Then

$$\hat{Y} = \underline{\Phi} \hat{\beta}_\lambda = \sum_{j=1}^{\text{rank}(\underline{\Phi})} v_j \downarrow \begin{matrix} <Y, v_j> \\ \frac{\sigma_j^2}{\sigma_j^2 + \lambda} \end{matrix} \quad \begin{matrix} \sigma_j \text{ } j^{\text{th}} \text{ singular} \\ \text{value of } \underline{\Phi} \end{matrix}$$

jth column of U

In contrast $\underline{\Phi} \hat{\beta} = \sum_{j=1}^{\text{rank}(\underline{\Phi})} v_j <Y, v_j>$

- Thm 3.7 The excess risk of $\hat{\beta}_\lambda$ is

$$\begin{aligned} \mathbb{E} [R(\hat{\beta}_\lambda)] - R(\beta^*) &= \sigma^2 \beta^* (\hat{\Sigma} + \lambda I_d)^{-2} \hat{\Sigma} \beta^* + \\ &\quad \underbrace{\frac{\sigma^2}{n} + \left(\sum_{j=1}^n \frac{\hat{\sigma}_j^2}{(\hat{\sigma}_j^2 + \lambda)} \right)}_{\text{jth eigenval. of } \hat{\Sigma}} \\ &= \text{Bias term} + \text{Variance term} \end{aligned}$$

- Bias is ↑ in λ variance is ↓ in λ.

- Task : choose optimal λ ! In Prop. 3.8 of

Bach's book you will see a convenient choice of λ

that minimizes an upper bound on the excess risk

$$\text{In particular with } \lambda = \frac{\sigma^2 \text{tr}(\hat{\Sigma})}{\|\beta^*\| \sqrt{n}}$$

$$\mathbb{E} [R(\hat{\beta}_\lambda)] - R(\beta^*) \leq \frac{\sigma^2 \text{tr}(\hat{\Sigma}) \|\beta^*\|_2}{\sqrt{n}}$$

- Remark : 1) This is a slow rate

- 2) we do not know σ or $\|\beta^*\|_2$

in) in practice you will do cross validation

- MIN(MAX
LOWER BOUND ON EXCESS RISK (Section 3.7 in Bach's book)