

Regression Notes

①

① Basics

Given $x_1, \dots, x_n \in \mathbb{R}^p$ & $y_1, \dots, y_n \in \mathbb{R}$, produce "good" β

Fitting = $\min_{\beta'} \|X\beta' - y\|^2$ $\min_{\beta} \|y - X\beta\|^2 + \lambda \|\beta\|^2$

"Ridge"

• Fitting data

• Prediction

• Parameter Estimation

$\beta_{OLS} = (X^T X)^{-1} X^T y$

$\beta_{\lambda} = (X^T X + \lambda I_p)^{-1} X^T y$

↑ Say: "assume invertible"

Fun Fact #1

Hat matrix $H = X(X^T X)^{-1} X^T$

Hat ["projection" on "influence"]

$\hat{y} = X\beta_{OLS}$
 $= X(X^T X)^{-1} X^T y$
 $= Hy$

$H_{ij} = x_i^T (X^T X)^{-1} x_j$

$H_{ii} \triangleq h_i$ the leverage of \bar{e}_i

$\sum_{i=1}^n h_i = \text{rank}$

Fun Fact #2

Influence of point \bar{e}_i .

Claim $h_i = \frac{\partial \hat{y}_i}{\partial y_i}$

Even stronger

$\beta_{-i} = (X^T X - x_i x_i^T)^{-1} (X^T y - x_i y_i)$

Sherman-Morrison formula

$= (X^T X)^{-1} + \frac{(X^T X)^{-1} x_i x_i^T (X^T X)^{-1}}{1 - x_i^T (X^T X)^{-1} x_i} (X^T y - x_i y_i)$

$= \beta_{OLS} + \frac{(X^T X)^{-1} x_i}{1 - h_i} (y_i - \langle x_i, \beta_{OLS} \rangle)$

$(A + uv^T)^{-1} = A^{-1} - \frac{A^{-1} u v^T A^{-1}}{1 + v^T A^{-1} u}$

$\langle x_i, \beta_{-i} \rangle = \langle x_i, \beta_{OLS} \rangle + \frac{H_{ii}}{1 - h_i} e_i$

Regression Notes

(2)

(2) Modeling Assumptions & Inference

Fixed design: X not random

$$y_i = \langle x_i, \beta \rangle + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$y = X\beta + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2 I)$$

Random design: $\{x_i\} \stackrel{iid}{\sim} \mathcal{N}(0, \Sigma)$

$$y = X\beta + \epsilon$$

Fixed design is easy:

Claim In fixed design setting, let $\tilde{\Sigma} = \frac{1}{n} X^T X$
 $\hat{\beta}_{OLS} = \beta + \mathcal{N}\left(0, \frac{\sigma^2}{n} \tilde{\Sigma}^{-1}\right)$

PF $\hat{\beta}_{OLS} = (X^T X)^{-1} X^T y$

$$= (X^T X)^{-1} X^T (X\beta + \epsilon)$$

$$= (X^T X)^{-1} X^T X \beta + (X^T X)^{-1} X^T \epsilon$$

$$= \beta + \mathcal{N}\left(0, (X^T X)^{-1} X^T X (X^T X)^{-1} \cdot \sigma^2\right)$$

$A \mathcal{N}(0, I) \stackrel{\Delta}{=} \mathcal{N}(0, AA^T)$

□

Regression Notes

(3)

Ask: everyone know about MLE & MAP?

For Gaussian data: $\operatorname{argmax}_{\beta'} P(y, X | \beta')$ ~~$\operatorname{argmax}_{\beta}$~~ $\xrightarrow{\text{indep of } \beta}$

$$\begin{aligned} &= \operatorname{argmax}_{\beta'} P(y | X, \beta') P(X | \beta') \\ &= \operatorname{argmax}_{\beta'} \mathcal{N}(X\beta', \sigma^2 I) \\ &= \operatorname{argmax}_{\beta'} \exp\left\{-\frac{1}{2\sigma^2} \|y - X\beta'\|^2\right\} \\ &= \operatorname{argmin} \|y - X\beta'\|^2 \end{aligned}$$

Maximum a posteriori $\operatorname{argmax}_{\beta'} P(\beta | y, X)$ $\operatorname{prior} P(\beta) = \mathcal{N}(0, \tau^2 I)$

$$\begin{aligned} &= \operatorname{argmax}_{\beta'} P(y | X, \beta') P(X | \beta') P(\beta') \\ &= \operatorname{argmax} \exp\left\{-\frac{1}{2\sigma^2} \|y - X\beta'\|^2 - \frac{1}{2\tau^2} \|\beta'\|^2\right\} \\ &= \operatorname{argmin} \|y - X\beta'\|^2 + \frac{\sigma^2}{\tau^2} \|\beta'\|^2 \end{aligned}$$

(3) Prediction (see: Hsu, Kakade, Zhang 11 & Mourouada, Rosasco 22)

In random design setting:

$$\begin{aligned} \mathbb{E}_{x, y} (\langle x, \hat{\beta} \rangle - y)^2 &= \mathbb{E} (\langle x, \hat{\beta} \rangle - \langle x, \beta \rangle - \eta)^2 \\ &= \mathbb{E} (\hat{\beta} - \beta)^T x x^T (\hat{\beta} - \beta) + \mathbb{E} \eta^2 \\ &= (\hat{\beta} - \beta)^T \Sigma (\hat{\beta} - \beta) + \sigma^2 \\ &= \|\Sigma^{1/2} (\hat{\beta} - \beta)\|^2 + \sigma^2 \end{aligned}$$

Regression Notes

(4)

Main Question What is ~~$\mathbb{E} \|\Sigma^{1/2}(\beta_n - \beta)\|^2$~~ ?

Hope: as in fixed design & OLS

$$\|\Sigma^{1/2}(\beta_{OLS} - \beta)\|^2 = \|\Sigma^{1/2} \cdot \mathcal{N}(0, \frac{\sigma^2 \hat{\Sigma}^{-1}}{n})\|^2$$

$$= \|\Sigma^{1/2} \frac{\sigma}{\sqrt{n}} \cdot \hat{\Sigma}^{1/2} \cdot \mathcal{N}(0, \mathbb{I})\|^2$$

$$\approx \frac{\sigma^2}{n} \cdot \|\mathcal{N}(0, \mathbb{I})\|^2 \approx \frac{\sigma^2 p}{n}$$

Assumptions on data-generating distribution

(A1) $\exists \beta \in \mathbb{R}^p$ s.t. $\mathbb{E}[y|x] = \langle \beta, x \rangle$ and $\text{Var}(y|x) \leq \sigma^2$

(A2) $\exists R$ s.t. $\|x\| \leq R$ almost surely

Thm [Montada & Rosasco, 22(?)]
Under A1 & A2, if $\lambda \approx R^2/n$

$$\mathbb{E} \left[\|\Sigma^{1/2}(\beta_n - \beta)\|^2 \right] \leq \lambda \|\beta\|^2 + \frac{\sigma^2 \text{Tr}[(\Sigma + \lambda \mathbb{I})^{-1} \Sigma]}{n}$$

More generally, $\forall \lambda > 0$

$$\mathbb{E} [J] \leq \left(1 + \frac{R^2}{\lambda n}\right)^2 \lambda \|\beta\|^2 + \left(1 + \frac{R^2}{\lambda n}\right) \frac{\sigma^2 \text{Tr}[(\Sigma + \lambda \mathbb{I})^{-1} \Sigma]}{n}$$

\hookrightarrow or $\inf_{\beta \in \mathbb{R}^p} \{ L(\beta) + \lambda \|\beta\|^2 - L(\hat{\beta}) \}$

Regression Notes

(5)

Lemma 1 Under $A1$ & $A2$,

$$\mathbb{E} \left[\|\Sigma^{1/2} (\beta_\lambda - \beta)\|^2 \right] \leq \lambda^2 \mathbb{E} \left[\langle (\hat{\Sigma} + \lambda \mathbb{I})^{-1} \Sigma (\hat{\Sigma} + \lambda \mathbb{I})^{-1} \beta, \beta \rangle \right]$$

bias

$$+ \frac{\sigma^2}{n} \mathbb{E} \left[\text{tr} \left((\hat{\Sigma} + \lambda \mathbb{I})^{-1} \Sigma \right) \right]$$

variance

Proof Idea
Key step

$$\beta_\lambda = (\hat{\Sigma} + \lambda \mathbb{I})^{-1} \frac{1}{n} \sum_{i=1}^n x_i y_i$$

$$= (\hat{\Sigma} + \lambda \mathbb{I})^{-1} \frac{1}{n} \sum_{i=1}^n x_i (x_i^T \beta + \eta_i)$$

$$= (\hat{\Sigma} + \lambda \mathbb{I})^{-1} \hat{\Sigma} \beta + (\hat{\Sigma} + \lambda \mathbb{I})^{-1} \frac{1}{n} \sum_{i=1}^n x_i \eta_i$$

$$\beta_\lambda - \beta = (\hat{\Sigma} + \lambda \mathbb{I})^{-1} \hat{\Sigma} \beta + (\hat{\Sigma} + \lambda \mathbb{I})^{-1} \frac{1}{n} \sum_{i=1}^n x_i \eta_i - (\hat{\Sigma} + \lambda \mathbb{I})^{-1} (\hat{\Sigma} + \lambda \mathbb{I}) \beta$$

$$= -\lambda (\hat{\Sigma} + \lambda \mathbb{I})^{-1} \beta + (\hat{\Sigma} + \lambda \mathbb{I})^{-1} \frac{1}{n} \sum_{i=1}^n x_i \eta_i$$

"□"

Lemma 2 Under $A2$,

$$\mathbb{E} \text{tr} \left[(\hat{\Sigma} + \lambda \mathbb{I})^{-1} \Sigma \right] \leq \left(1 + \frac{R^2}{\lambda n} \right) \text{tr} \left[(\hat{\Sigma} + \lambda \mathbb{I})^{-1} \Sigma \right]$$

PF $x_{n+1} \sim \mathcal{N}(0, \Sigma)$, indep of iid from x_1, \dots, x_n

$$\mathbb{E}_{x_1, \dots, x_n} \text{tr} \left[(\hat{\Sigma} + \lambda \mathbb{I})^{-1} \Sigma \right] = \mathbb{E}_{x_1, \dots, x_n, x_{n+1}} \text{tr} \left[(\hat{\Sigma} + \lambda \mathbb{I})^{-1} x_{n+1} x_{n+1}^T \right]$$

$$= n \mathbb{E}_{x_1, \dots, x_n, x_{n+1}} \left[\langle (n \hat{\Sigma} + \lambda n \mathbb{I})^{-1} x_{n+1}, x_{n+1} \rangle \right]$$

⇒

Sherman-Morrison again

$$\begin{aligned}
 & x_{n+1}^T \left(n \hat{\Sigma} + x_{n+1} x_{n+1}^T + \lambda n I \right)^{-1} x_{n+1} \\
 &= x_{n+1}^T \left(n \hat{\Sigma} + \lambda n I \right)^{-1} x_{n+1} - \frac{x_{n+1}^T \left(n \hat{\Sigma} + \lambda n I \right)^{-1} x_{n+1} x_{n+1}^T \left(n \hat{\Sigma} + \lambda n I \right)^{-1} x_{n+1}}{1 + x_{n+1}^T \left(n \hat{\Sigma} + \lambda n I \right)^{-1} x_{n+1}} \\
 &= x_{n+1}^T \left(n \hat{\Sigma} + \lambda n I \right)^{-1} x_{n+1} \left(1 - \frac{x_{n+1}^T \left(n \hat{\Sigma} + \lambda n I \right)^{-1} x_{n+1}}{1 + x_{n+1}^T \left(n \hat{\Sigma} + \lambda n I \right)^{-1} x_{n+1}} \right)
 \end{aligned}$$

So

$$n \mathbb{E} \left[x_{n+1}^T \left(n \hat{\Sigma} + \lambda n I \right)^{-1} x_{n+1} \right]$$

x_1, \dots, x_n
 x_{n+1}

$$\mathbb{E} \left(\begin{array}{c} 1 - \frac{\mathbb{E}}{1+\mathbb{E}} \\ \frac{1+\mathbb{E} - \mathbb{E}}{1+\mathbb{E}} \end{array} \right)$$

$$\begin{aligned}
 &= n \mathbb{E} \left[\left(1 + x_{n+1}^T \left(n \hat{\Sigma} + \lambda n I \right)^{-1} x_{n+1} \right)^{-1} x_{n+1}^T \left((n+1) \hat{\Sigma}_{n+1} + \lambda n I \right)^{-1} x_{n+1} \right] \\
 &\quad \hookrightarrow \leq \left(1 + \frac{R^2}{\lambda n} \right)
 \end{aligned}$$

$$\leq n \left(1 + \frac{R^2}{\lambda n} \right) \mathbb{E} \left[x_{n+1}^T \left((n+1) \hat{\Sigma}_{n+1} + \lambda n I \right)^{-1} x_{n+1} \right]$$

$$\stackrel{\text{key step}}{=} n \left(1 + \frac{R^2}{\lambda n} \right) \frac{1}{n+1} \sum_{i=1}^{n+1} \mathbb{E} \left[x_i^T \left((n+1) \hat{\Sigma}_{n+1} + \lambda n I \right)^{-1} x_i \right]$$

$$= n \left(1 + \frac{R^2}{\lambda n} \right) \frac{1}{n+1} \sum_{i=1}^{n+1} \mathbb{E} \left[\text{tr} \left[\left((n+1) \hat{\Sigma}_{n+1} + \lambda n I \right)^{-1} x_i x_i^T \right] \right]$$

$$= n \left(1 + \frac{R^2}{\lambda n} \right) \mathbb{E} \left[\text{tr} \left[\left((n+1) \hat{\Sigma}_{n+1} + \lambda n I \right)^{-1} \hat{\Sigma}_{n+1} \right] \right]$$

$\mathbb{E} \mathbb{1}$

$\frac{1}{n} \lambda I$

Observe:
 $\hat{\Sigma}_{n+1} \preceq \hat{\Sigma}_{n+1}$

Regression Notes

6

$$\left(\left(1 + \frac{1}{n} \right) \hat{\Sigma}_{n+1} + \lambda \mathbb{I} \right)^{-1} \left\{ \left(\hat{\Sigma}_{n+1} + \lambda \mathbb{I} \right)^{-1} \right.$$

Claim: $f: A \mapsto \text{tr} \left[(A + \lambda \mathbb{I})^{-1} A \right]$ is concave.

Idea Pf idea: $\text{tr} \left[(A + \lambda \mathbb{I})^{-1} A \right] = \sum_{i=1}^p \frac{\lambda_i}{\lambda_i + \lambda}$

