

3) Regularizacija linearnog modela [ridge regresija]

Notacija

[prvo želimo pokazati da LS procjena može biti loša i kada je Lin. model tačan !?]

• vektore u \mathbb{R}^k ($k \in \mathbb{N}$) shvaćamo kao vektor-stupce

$$X = \begin{bmatrix} \text{---} x^{(1)\top} \text{---} \\ \vdots \\ \text{---} x^{(n)\top} \text{---} \end{bmatrix} = \begin{bmatrix} x_{11}, \dots, x_{1p} \\ \vdots \\ x_{n1}, \dots, x_{np} \end{bmatrix} \in \mathbb{R}^{n \times p}$$

matrica dizajna

• $y = (y_1, \dots, y_n)^\top \in \mathbb{R}^{n \times 1}$

• $x_j = (x_{1j}, \dots, x_{nj})^\top \in \mathbb{R}^{n \times 1}$, $j = 1, \dots, p$

↳ j -ti stupac od X (mjesto: koefficijenti x_j u svim $x^{(i)}$, $i = 1, \dots, n$)

(često u X dodajemo $x_0 = (1, \dots, 1)^\top \in \mathbb{R}^{n \times 1}$
 x_{10}, \dots, x_{n0})

Lin. model: biramo $\hat{f} = \hat{f}(x)$ iz $\mathcal{F} = \{f_\beta : \beta \in \Theta\}$ uz $\Theta \subseteq \mathbb{R}^{p+1}$;

$$f_\beta(x) = X^\top \cdot \beta, \quad X = (x_0, \dots, x_p) \in \mathbb{R}^{p+1}. \quad (3.1)$$

[dopušteno: (i) $x_2 = f(x_1)$, (ii) interakcije $x_3 = x_1 \cdot x_2$, (iii) kvalitativne koefficijente.]

1) Metoda najmanjih kvadrata → LS ["least squares"]

$$\hat{\beta}^{LS} = (\hat{\beta}_0^{LS}, \dots, \hat{\beta}_p^{LS}) := \arg \min_{\beta \in \mathbb{R}^{p+1}} \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - f_\beta(x^{(i)}))^2}_{L_T(f_\beta) \text{ uz } L = (y - \hat{y})^2} \quad (3.2)$$

$$= \arg \min_{\beta} \|y - X\beta\|_2^2$$

Pretpostavimo da su $x_0, \dots, x_p \in \mathbb{R}^m$ lin. nezavisni, tj. da je $\boxed{24}$
 $X^T X \in \mathbb{R}^{(p+1) \times (p+1)}$ poz. definitna (simetrična) matrica.
 (nužno je $\boxed{m \geq p+1}$)

• $\hat{\beta}^{OLS} = (X^T X)^{-1} X^T y = \frac{1}{n} \hat{\Sigma}^{-1} X^T y$, (3.3)

gdje je $\boxed{\hat{\Sigma} := \frac{1}{n} X^T X}$. [kada su stupci centrirani, ovo je točno varijakozna kovarijacijska matrica]

• $\hat{y} := (\hat{\beta}_{OLS}(x^{(1)}), \dots, \hat{\beta}_{OLS}(x^{(m)}))^T = X \cdot \hat{\beta}^{OLS} = H y$, (3.4)

uz $H = X (X^T X)^{-1} X^T$ ["hat matrix"],

te je \hat{y} ort. projekcija vektora y na $\text{span}(x_0, \dots, x_p)$.

[Nap. U slučaju da x_0, \dots, x_p nisu lin. nezavisni (npr. $m < p+1$ ili $x_j = a \cdot x_i$) $\hat{\beta}^{OLS}$ nije jedinstven. Ipak, \hat{y} je i dalje ortogonalna projekcija od y na $\text{span}(x_0, \dots, x_p)$.]

3.2 Statistička svojstva od $\hat{\beta}^{OLS}$

[minimalni pretpostaviti model za (X, y)]

Pretpostavimo da $\boxed{\exists \beta_* \in \mathbb{R}^{p+1}}$ t.d.

$y_i = (X^{(i)})^T \beta_* + \epsilon_i, i=1, \dots, n$ (3.5)

pri čemu su

[radi jednostavnosti]

• $X^{(i)}, i=1, \dots, n$ neslučajni (tzv. "fiksni dizajn")

• $\epsilon_i, i=1, \dots, n$, nezavisni t.d. $\boxed{E[\epsilon_i] = 0}$, $\boxed{\text{Var}(\epsilon_i) = \sigma^2}$, $\forall i=1, \dots, n$.

za neko $\sigma^2 > 0$.

Napz-1 u (3.5) pretpostavljamo da je lin. model tačan, tj. da nema greške aproksimacije ("bias") 25

(3.5) možemo zapisati kao [tikovs]

$$Y := \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = X \beta_* + \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix} =: X \beta_* + \epsilon, \quad (3.6)$$

uz $E[\epsilon] = 0, \text{Var}(\epsilon) = \sigma^2 I_n$.

$\Rightarrow \hat{\beta}^{OLS} = \beta_* + \frac{1}{n} \sum^{-1} X^T \epsilon$ jedini slučajno

(3.3) Rechno

- $E[\hat{\beta}^{OLS}] = \beta_*$
- $\text{Var}(\hat{\beta}^{OLS}) = \frac{\sigma^2}{n} \sum^{-1}$

(3.7)

Tipično

$$\hat{\sigma}^2 := \frac{1}{n-p-1} \|Y - \hat{Y}\|_2^2 = \frac{RSS(\hat{\beta})}{n-p-1} \quad (3.8)$$

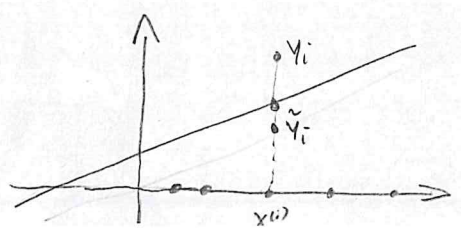
3.2.1 Testna greška \rightarrow F. Brech, (3.5)

Neka je $\tilde{Y} := X \beta_* + \tilde{\epsilon}$, pri čemu su $\epsilon_i \tilde{\epsilon}$ (tj. $Y_i \tilde{Y}_i$) nezavisni, testni skup.

Promatramo tzv. testnu grešku unutar uzorka ("in-sample")

$$L^{in}(\beta) := E \left[\frac{1}{n} \sum_{i=1}^n (\tilde{Y}_i - \beta(x^{(i)}))^2 \right] = E \left[\frac{1}{n} \| \tilde{Y} - X \beta \|_2^2 \right], \quad (3.9)$$

za $\forall \beta \in \mathbb{R}^{p+1}$.



Prp. 3.1 | Ako je $\hat{\beta} = \hat{\beta}(Y, X)$ procjenitelj za β_0 ,
(BVT)

$$E[L^m(\hat{\beta})] = \sigma^2 + \overbrace{\|E[\hat{\beta}] - \beta_0\|_{\Sigma}^2}^{\text{"pristranost"}} + \overbrace{E[\|\hat{\beta} - E[\hat{\beta}]\|_{\Sigma}^2]}^{\text{"varijanca"}}, \quad (3.10)$$

gdje $E = E_Y$ te $\|\beta\|_{\Sigma}^2 := \beta^T \hat{\Sigma} \beta, \beta \in \mathbb{R}^{p+1}$.

Dokaz 1

(1) Iz (3.6) imamo, $\forall \beta \in \mathbb{R}^{p+1}$,
za \tilde{y}

[uzimamo u obzir distribuciju podataka]

$$\begin{aligned} m \cdot L^m(\hat{\beta}) &= E[\|\tilde{y} - X\beta\|_2^2] = E[\|X(\beta_0 - \beta) + \tilde{\epsilon}\|_2^2] \\ &= E[\|a+b\|_2^2 = \langle a+b, a+b \rangle = \|a\|_2^2 + \|b\|_2^2 + 2 \cdot a^T b, a, b \in \mathbb{R}^{p+1}] \\ &= E[\|X(\beta_0 - \beta)\|_2^2 + \|\tilde{\epsilon}\|_2^2 + \underbrace{2(X(\beta_0 - \beta))^T \tilde{\epsilon}}_{\text{mije slučajno}}] \\ &= [E[\|\tilde{\epsilon}\|_2^2] = m \cdot \sigma^2, E[\tilde{\epsilon}] = 0 \in \mathbb{R}^n] \\ &= m \cdot \sigma^2 + (\beta_0 - \beta)^T X^T X (\beta_0 - \beta) \end{aligned}$$

min se postiže za $\beta_k i = \hat{\beta}_k$

$$\Rightarrow L^m(\hat{\beta}) = \sigma^2 + \|\beta_0 - \hat{\beta}\|_{\Sigma}^2 \quad (3.11)$$

(2)

$$E[L^m(\hat{\beta})] \stackrel{(3.11)}{=} \sigma^2 + E[\|\beta_0 - \hat{\beta}\|_{\Sigma}^2]$$

$$\begin{aligned} E[\|\beta_0 - \hat{\beta}\|_{\Sigma}^2] &= E[\|(\beta_0 - E[\hat{\beta}]) + (E[\hat{\beta}] - \hat{\beta})\|_{\Sigma}^2] \\ &= E[\|a+b\|_{\Sigma}^2 \stackrel{(3.11)}{=} \|a\|_{\Sigma}^2 + \|b\|_{\Sigma}^2 + 2 a^T \hat{\Sigma} b] \\ &= E[\|\beta_0 - E[\hat{\beta}]\|_{\Sigma}^2] + E[\|E[\hat{\beta}] - \hat{\beta}\|_{\Sigma}^2] + \end{aligned}$$

$$+ 2 \underbrace{E[(\beta_0 - E\hat{\beta})^T \hat{\Sigma} (\hat{E}\hat{\beta} - \hat{\beta})]}_{\text{nije slučajno}}$$

[2:

$$= (\beta_0 - E\hat{\beta})^T \hat{\Sigma} \underbrace{E[\hat{E}\hat{\beta} - \hat{\beta}]}_{= 0}$$

Prop. 3.2 (!)

$$E[L^{(n)}(\hat{\beta}_{OLS})] = \sigma^2 + 0 + \boxed{\frac{\sigma^2(p+1)}{n}} \quad (3.12)$$

Dokaz. | Sljedi iz Prop. 3.1 jer $E[\hat{\beta}_{OLS}] = \beta_0$ te

$$E[\|\hat{\beta}_{OLS} - E[\hat{\beta}_{OLS}]\|_{\hat{\Sigma}}^2] = E[\underbrace{(\hat{\beta}_{OLS} - E[\hat{\beta}_{OLS}])^T \hat{\Sigma} (\hat{\beta}_{OLS} - E[\hat{\beta}_{OLS}])}_{\in \mathbb{R}}]$$

$$= E[\text{tr}(\cdot)] = E[\text{tr}((\hat{\beta}_{OLS} - E[\hat{\beta}_{OLS}]) (\hat{\beta}_{OLS} - E[\hat{\beta}_{OLS}])^T \hat{\Sigma})]$$

nije slučajno

$$= \text{tr}(\underbrace{\text{Var}(\hat{\beta}_{OLS})}_{\substack{= \frac{\sigma^2}{n} \hat{\Sigma}^{-1} \\ (3.7)}} \cdot \hat{\Sigma}) = \frac{\sigma^2}{n} \text{tr}(\mathbb{I}_{p+1}) = \frac{\sigma^2(p+1)}{n}$$

Alternativ,

$$E[\|\hat{\beta}_{OLS} - E[\hat{\beta}_{OLS}]\|_{\hat{\Sigma}}^2] = E[\|\underbrace{(\hat{\beta}_{OLS} - \beta_0)}_{= \text{tr}(\epsilon^T H \epsilon)}\|_{\hat{\Sigma}}^2] = E[\| (X^T X)^{-1} X^T \epsilon \|_{\hat{\Sigma}}^2]$$

$$= \dots = \frac{1}{n} E[\epsilon^T H \epsilon] = \frac{1}{n} E[\text{tr}(\epsilon \epsilon^T H)] = \frac{\sigma^2}{n} \text{tr}(H) = \frac{\sigma^2}{n} (p+1)$$

$E[\epsilon \epsilon^T] = \sigma^2 \mathbb{I}_n$
 $\downarrow + \text{tr}(H) = p+1$

(i) Iz $\hat{Y} = HY$ te Prop. 2.7 sledi da

$$df(LS) = \text{tr}(H) = p + 1 = \# \text{ procenjenih parametara.}$$

H ort. projektor na span $\{x_0, \dots, x_p\}$
lin. nez.

(ii) Nadalje, iz Prop (3.1) i (3.2) imamo $(E[L_T(\hat{\beta})] = E_Y[L^*(\hat{\beta})])$,

$$E[L_T(\hat{\beta}_{LS})] = \sigma^2 \left(1 - \frac{\sigma^2(p+1)}{n} \right)$$

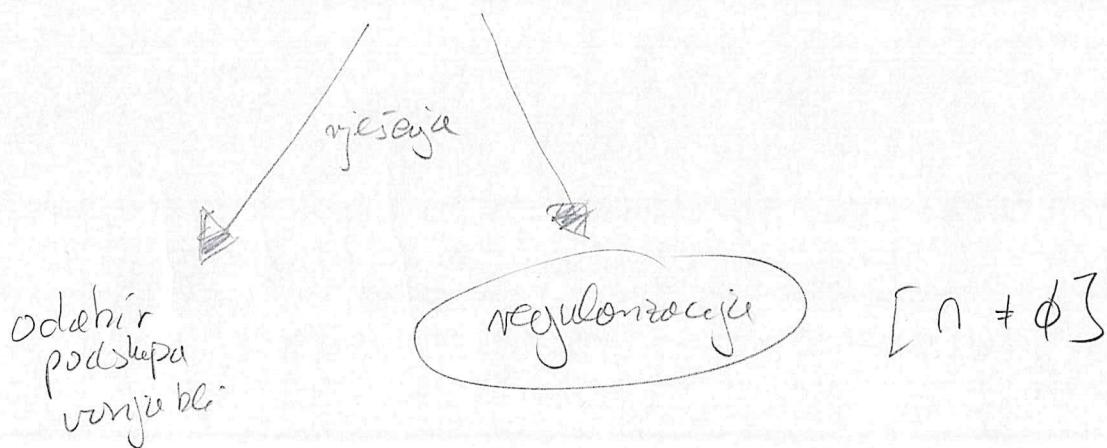
[training greška]

$$= \sigma^2 \frac{n-p-1}{n} \left[\xrightarrow{p+1 \rightarrow n} 0 \right]$$

[Specijalno, za $\hat{\sigma}^2$ iz (3.2) imamo $E[\hat{\sigma}^2] = \sigma^2$.]

(iii) Dakle, čak i kad je lin. model točan (pa je pristranost = 0), varijanca [a somim time i ukupna greška] je velika [u odnosu na σ^2] ako je $\left(\frac{p+1}{n}\right)$ velik!

[čak i da svih p koeficijenta snažno utječu na odziv!]



(3.3) Ridge regression

Pr (Motivacija) [Zasto ne uvijek naprostranstvo ?]

Pretp. da z_1, \dots, z_n njd t.d. $z_i \sim N(\theta, \sigma^2)$.

• $MSE(\bar{z}) := E[(\bar{z} - \theta)^2] = \text{Var}(\bar{z}) + 0 = \frac{\sigma^2}{n} + 0$

• $MSE(\bar{z}/2) := E[(\frac{\bar{z}}{2} - \theta)^2] = \text{Var}(\bar{z}/2) + (E[\bar{z}/2] - \theta)^2 = \frac{\sigma^2}{4n} + \frac{\theta^2}{4}$

$\hookrightarrow MSE(\bar{z}) \leq MSE(\bar{z}/2)$

$n \geq \frac{3\sigma^2}{\theta^2}$ velike obr $\theta \approx 0$.

Za $\lambda \geq 0$ defem.

" L_2 regularizacija"

$\hat{\beta}_{\lambda}^r = \hat{\beta}^r = (\hat{\beta}_0^r, \dots, \hat{\beta}_p^r) = \underset{\beta \in \mathbb{R}^{p+1}}{\text{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^m (y_i - f_{\beta}(x_i^{(m)}))^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$

$= \underset{\beta \in \mathbb{R}^{p+1}}{\text{argmin}} \left\{ \text{RSS}(f_{\beta}) + m \lambda \sum_{j=1}^p \beta_j^2 \right\}$.

(3.13)

ne uključuje β_0 !

3.4 | $\forall \lambda \geq 0, \forall \beta_1, \dots, \beta_p \in \mathbb{R}$,

$\underset{\beta_0 \in \mathbb{R}}{\text{argmin}} \left\{ \sum_{i=1}^m (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + m \lambda \sum_{j=1}^p \beta_j^2 \right\} = \bar{y} - \sum_{j=1}^p \beta_j \bar{x}_j$,

je je $\bar{y} = \frac{1}{m} \sum_{i=1}^m y_i, \bar{x}_j = \frac{1}{m} \sum_{i=1}^m x_{ij}, j=1, \dots, p$.

$(\hat{\beta}_1^r, \dots, \hat{\beta}_p^r) = \underset{\beta \in \mathbb{R}^p}{\text{argmin}} \left\{ \sum_{i=1}^m (y_i - \bar{y} - \sum_{j=1}^p (x_{ij} - \bar{x}_j) \beta_j)^2 + m \lambda \sum_{j=1}^p \beta_j^2 \right\}$

Zbog Prop. 3.4 u nastavku pretpostavljamo da je

$$\bar{x}_j = 0, \quad j = 1, \dots, p, \quad \bar{y} = 0, \quad (3.14)$$

te rješavamo

$$\hat{\beta}^r = (\hat{\beta}_1^r, \dots, \hat{\beta}_p^r) = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \|y - X\beta\|_2^2 + n\lambda \|\beta\|_2^2 \right\} \quad (3.15)$$

(vidi (3.24))

Prop. 3.5 | Uz centriranje, kovarijate tipično još i skaliramo t.d.

imaju varijancu 1 jer je kazna $\|\beta\|_2^2$

"nepošteni" ako kovarijate nisu u istim mjerim jedinicama.

Prop 3.6 | $\forall \lambda \geq 0$

$$\left[\left(\frac{y - \bar{y}}{\Delta y} \right) = \hat{\beta}_1^r \cdot \left(\frac{x_1 - \bar{x}_1}{\Delta x_1} \right) + \dots + \hat{\beta}_p^r \cdot \left(\frac{x_p - \bar{x}_p}{\Delta x_p} \right) \Rightarrow y = \hat{\beta}_0^r + \hat{\beta}_1^r \frac{\Delta y}{\Delta x_1} x_1 + \dots + \hat{\beta}_p^r \frac{\Delta y}{\Delta x_p} x_p \right]$$

$$\hat{\beta}^r = \frac{1}{n} (\hat{\Sigma} + \lambda I)^{-1} X^T y, \quad (3.16)$$

Dokaz | (Dz) \rightarrow analogno kao za $\hat{\beta}^{OLS}$.

Prop. | (i) za $\lambda > 0$, $\hat{\beta}^r$ je uvijek dobro definiran (dobro, i u

slučaju $p > n$!)

- (ii) • $\hat{\beta}^r = \hat{\beta}^{OLS}$ za $\lambda = 0$,
 - $\hat{\beta}^r \xrightarrow{\lambda \rightarrow \infty} 0$.
- } hiperparameter λ kontrolira efekt smanjenja (engl. shrinkage) koeficijenta prema 0.

ridge - primjer - smanjenja - koef

Prop. 3.7 | Pod pretpostavkama u (3.5), $\forall \lambda \geq 0$,

$$E \left[L^{(n)}(\hat{\beta}_{\lambda}^r) \right] = d^2 + \lambda^2 \beta_*^T (\hat{\Sigma} + \lambda I_p)^{-2} \hat{\Sigma} \beta_* + \frac{d^2}{n} \operatorname{tr} \left[\hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2} \right]. \quad (3.17)$$

varijanca

$$\hat{\beta}_\lambda^r = \frac{1}{n} (\hat{\Sigma} + \lambda I_p)^{-1} X^T Y + \frac{1}{n} (\hat{\Sigma} + \lambda I_p)^{-1} X^T \varepsilon \quad (a)$$

$$Y = X\beta_\alpha + \varepsilon$$

iz Prop. 3.1, $\mathbb{E}[\| \hat{\beta}_\lambda^r \|^2] - d^2 = B + V$, gdje

$$B = \mathbb{E}[\| \hat{\beta}_\lambda^r \|^2] - \beta_\alpha^T \beta_\alpha$$

$$\begin{aligned} \mathbb{E}[\hat{\beta}_\lambda^r] &= \frac{1}{n} (\hat{\Sigma} + \lambda I_p)^{-1} X^T X \beta_\alpha = (\hat{\Sigma} + \lambda I_p)^{-1} [\hat{\Sigma} \beta_\alpha + \lambda I_p \beta_\alpha] \\ &= \beta_\alpha - \lambda (\hat{\Sigma} + \lambda I_p)^{-1} \beta_\alpha \end{aligned}$$

$\mathbb{E}[\varepsilon] = 0$

$$\begin{aligned} \Rightarrow B &= \mathbb{E}[\| \lambda (\hat{\Sigma} + \lambda I_p)^{-1} \beta_\alpha \|^2] \\ &= \lambda^2 \beta_\alpha^T (\hat{\Sigma} + \lambda I_p)^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda I_p)^{-1} \beta_\alpha \\ &= \lambda^2 \beta_\alpha^T (\hat{\Sigma} + \lambda I_p)^{-2} \hat{\Sigma} \beta_\alpha \end{aligned}$$

$\hat{\Sigma}^T = \hat{\Sigma}$ $\hat{\Sigma} = (\hat{\Sigma} + \lambda I_p)^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda I_p)$

$$V = \mathbb{E}[\| \hat{\beta}_\lambda^r - \mathbb{E}[\hat{\beta}_\lambda^r] \|^2]$$

$$= \mathbb{E}[\| \frac{1}{n} (\hat{\Sigma} + \lambda I_p)^{-1} X^T \varepsilon \|^2]$$

$$= \mathbb{E}[\text{tr}(\frac{1}{n} \varepsilon^T X^T (\hat{\Sigma} + \lambda I_p)^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda I_p)^{-1} X^T \varepsilon)]$$

$$= \frac{d^2}{n} \text{tr}(\hat{\Sigma} (\hat{\Sigma} + \lambda I_p)^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda I_p)^{-1})$$

$$\mathbb{E}[\varepsilon \varepsilon^T] = d^2 I_n$$

now done

$$= \frac{d^2}{n} \text{tr}(\hat{\Sigma}^2 (\hat{\Sigma}^2 + \lambda I_p)^{-2})$$

Uopz-1 (i) Pokazat ćemo da je

$$\text{tr} \left(\hat{\Sigma}^{-2} (\hat{\Sigma} + \lambda I)^{-2} \right) = \sum_{j=1}^p \frac{\lambda_j^2}{(\lambda_j + \lambda)^2}$$

gde su $\lambda_1, \dots, \lambda_p \geq 0$ svoj. vrijednosti od $\hat{\Sigma}$.

\Rightarrow
 $\forall \lambda > 0$

varijanca $< \frac{d^2}{n} \cdot P$, te
varijanca $\searrow 0$, za $\lambda \rightarrow \infty$.

[vidi (3.11)]

(ii) λ -norma
norma = $\beta_\alpha^\top (\frac{1}{\lambda} I + I_p)^{-2} \hat{\Sigma} \beta_\alpha \nearrow \beta_\alpha^\top \hat{\Sigma} \beta_\alpha = \|\beta_\alpha\|_{\hat{\Sigma}}^2$
 $= \|\beta_\alpha - 0\|_{\hat{\Sigma}}^2$
($\neq 0!$) kada $\lambda \rightarrow \infty$. \square

(iii) Može se pokazati da postoji $\lambda_\alpha > 0$ t.d.

$$E [L^{in}(\lambda_{\beta_\alpha}^{\hat{\Sigma}})] \searrow E [L^{in}(\lambda_{\beta_\alpha}^{\hat{\Sigma}_r})]$$

\hookrightarrow u praksi, λ odabiremo CV metodom.

[λ_α ne možemo eksplicitno računati jer ovisi o nepoznatim veličinama.]



3.3.1 Kako ridje smanjuje procjene \rightarrow rekonstruirati skraćeno 133

Prep. da je

$$\widehat{\text{Cov}}(x_j, x_k) = \frac{1}{n} x_j^T x_k = 0, \quad \forall j \neq k$$

Zbog $\bar{x}_j = 0, \forall j,$

$$\widehat{\Sigma} = \frac{1}{n} X^T X = \text{variančna kovarijacijska matrica}$$

od $x^{(1)}, \dots, x^{(n)}$

te

$$\widehat{\Sigma} = \text{diag} \left(\underbrace{\widehat{\text{Var}}(x_1)}_{\frac{\|x_1\|_2^2}{n}}, \dots, \widehat{\text{Var}}(x_p) \right) =: \text{diag}(\lambda_1, \dots, \lambda_p)$$

Sada imamo

[isto kao kod mjerenja y samo na x_j !]

$$\hat{\beta}_j^{ls} = \frac{x_j^T y}{\|x_j\|_2^2}, \quad j=1, \dots, p$$

(3.18)

$$\hat{y}_{ols} = \text{Proj}_{x_1, \dots, x_p}(y)$$

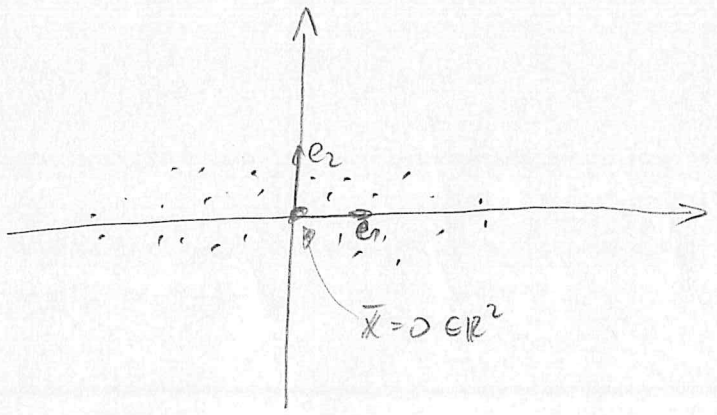
$$\hat{\beta}_j^{rr} = \frac{\|x_j\|_2^2}{\|x_j\|_2^2 + m\lambda} \cdot \hat{\beta}_j^{ls} = \frac{\lambda_j}{\lambda_j + \lambda} \cdot \hat{\beta}_j^{ls}, \quad j=1, \dots, p$$

(3.16) DT

≤ 1 za $\lambda > 0$, te

$$\frac{\lambda_j}{\lambda_j + \lambda} \rightarrow 0 \text{ kad } \lambda_j = \widehat{\text{Var}}(x_j) \rightarrow 0$$

prim. | $p=2, \text{Var}(x_1) \gg \text{Var}(x_2), m$ nije velik



Intuitivno, ako prilagodavamo lin. ploku (y -ova je treća koordinata), $\text{Var}(\hat{\beta}^{ls})$ bit će velika je procjena negitna u smjeru e_2 jer ovisi o T

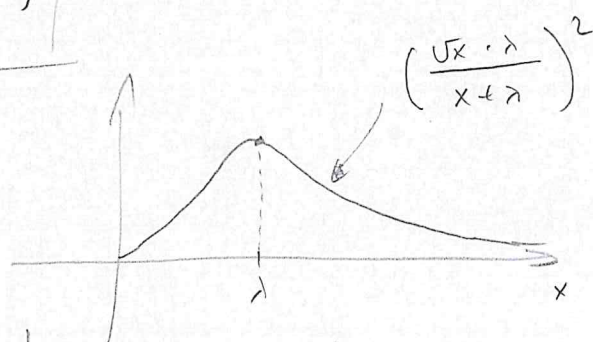
• Ridge smanjuje tu varijancu teku da smanjuje ("shrink") | 34

procjene, i to više u smanjenju e_2 .

→ dobivamo prostornost čija veličina ovisi o tome je li stvarni nagib u smanjenju e_2 (tj. β_2^*), velik li me!

• Iz (3.17) dobivamo da je

$$\text{bias} \underset{\substack{\uparrow \\ \text{DZ}}}{=} \sum_{j=1}^p \left(\frac{\sqrt{s_j \cdot \lambda}}{s_j + \lambda} \right)^2 \cdot (\beta_j^*)^2$$



• dakle veliku prostornost ako β_j^* velik te s_j nije prevelika (niti prevelika)

?

→ ridge - simulacijski - primjer . pdf

3.3.2) Kako vidje smanjuje koeficijente? - općeniti slučaj - [vidi 6. stranu]

SVD od $X \in \mathbb{R}^{n \times p}$ je

$$X = U \cdot D \cdot V^T, \quad (\text{SVD})$$

$n \times p \quad n \times p \quad p \times p \quad p \times p$

pri čemu

- $U^T U = I_p$
 - $V V^T = V^T V = I_p$
- stepci $u_1, \dots, u_p \in \mathbb{R}^{n \times 1}$ od U i $v_1, \dots, v_p \in \mathbb{R}^{p \times 1}$ su ortogonalni.

- $D = \text{diag}(d_1, \dots, d_p)$ sa $d_1 \geq d_2 \geq \dots \geq d_p \geq 0$
["singularne vrijednosti"]

$$\hat{\Sigma} = \frac{1}{n} X^T X = V \frac{D^2}{n} V^T$$

(SVD)

tj. $\lambda_j = \frac{d_j^2}{n}, j=1, \dots, p$ su svojstv. vrijednosti od $\hat{\Sigma}$, a

$v_j, j=1, \dots, p$ odgovarajuć. svojstv. vektori.

tzv. svojst. glavnih komponenti (od X)

↳ pref. da $\text{rang}(X) = p$ (te $n > p$) $\Rightarrow (d_p > 0)$.

$$Z_j := X v_j = \begin{bmatrix} v_j^T x^{(1)} \\ \vdots \\ v_j^T x^{(n)} \end{bmatrix} \in \mathbb{R}^{n \times 1}, j=1, \dots, P \quad (3.15)$$

Zovemo glavne komponente (od X).

$$(\text{Proj}_{v_j}(x) = \underline{v_j^T \cdot x} \cdot v_j, \text{ za } x \in \mathbb{R}^p)$$

Pop-1 vjedi, $\forall j=1, \dots, P$,

$$v_j = \underset{\substack{v \in \mathbb{R}^p \\ \|v\|_2=1}}{\text{arg max}} \widehat{\text{Var}}(X v),$$

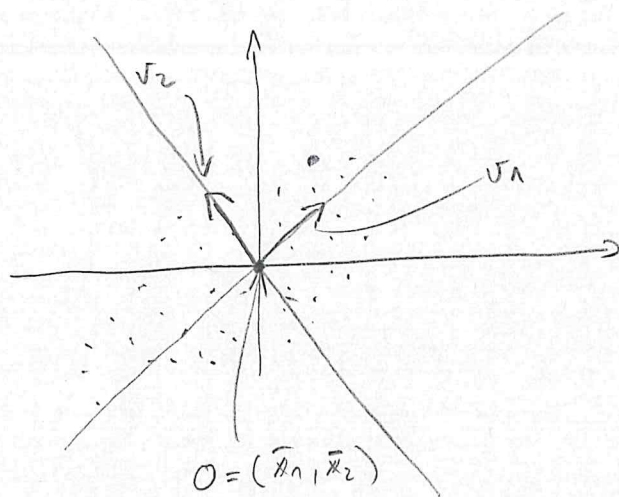
↖ varianca varijansa od $v^T x^{(i)}, i=1, \dots, n$

uz uvjet

$$v_j^T v_l = 0, \forall l=1, \dots, j-1.$$

[bez dokaza → Primjenjena statistika]

mpm. $p=2$



v_1 je smer v koji maksimizira varijancu projekcija od $x^{(1)}, \dots, x^{(n)}$ na v , v_2 je smer $\{v \perp v_1\}$ koji maksimizira \square
 ————, itd.

Neka je

$$Z := X V = \begin{bmatrix} | & & | \\ z_1 & \dots & z_p \\ | & & | \end{bmatrix} \in \mathbb{R}^{n \times p}, \quad (2.20)$$

z_j je podvektor u bazi u_1, \dots, u_p .

(SVD) $\Rightarrow Z = U \cdot D, z_j$.

$$z_j = d_j \cdot u_j, j=1, \dots, p \quad (2.21)$$

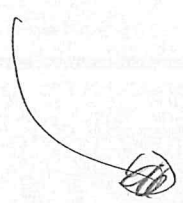
Specijalno,

$$[z_j] = \begin{pmatrix} \bar{x}_j \\ \vdots \\ 0 \end{pmatrix} \cdot u_j = 0, u_j$$

$$\|z_j\|^2 = z_j^T z_j = d_j^2, u_j, \quad (2.22)$$

$$\hat{\text{Var}}(z_j) = \frac{d_j^2}{n} = \underline{\underline{\sigma_j^2}}, u_j$$

$$\hat{\text{Cov}}(z_j, z_h) = \frac{1}{n} z_j^T z_h = 0, u_j \neq u_h \quad (2.23)$$



Ridge i glavne komponente

$\forall \beta \in \mathbb{R}^p,$

$X = UDV^T = ZV^T$

[38]

$\|y - X\beta\|_2^2 + \lambda n \|\beta\|_2^2 = \|y - Z \cdot \alpha\| + \lambda n \|\alpha\|_2^2$

gdje je $\alpha = V^T \beta$ ($V^2V = I \Rightarrow \|\alpha\|_2 = \|\beta\|_2$).



$\frac{1}{n} Z^T Z = \text{diag}(\lambda_1, \dots, \lambda_p)$
(2,22), (3,23)

[rekurezivni slučaj]

(3.18)

$\hat{\alpha}_j^r = \frac{\lambda_j}{\lambda_j + \lambda} \cdot \hat{\alpha}_j^{ls}, j = 1, \dots, p$

$(\hat{\beta}^r = V \cdot \hat{\alpha}^r)$

Općenito, vidje smo da je ls projekcije na v_1, \dots, v_p , te je smjernice veće za $\hat{\alpha}_j^{ls}$ za koje je $\widehat{Var}(\hat{\alpha}_j) = \lambda_j$ manje.

Uop. 3.9 | $d(\text{ridge}) = ?$

Imamo $\forall \lambda \geq 0,$

(3.16) $\text{diag}(\lambda_1, \dots, \lambda_p)$

$\hat{y}_\lambda^r = X \hat{\beta}_\lambda^r = Z \cdot \hat{\alpha}_\lambda^r = Z (Z^T Z + n\lambda I_p)^{-1} Z^T y$

$= U \cdot \text{diag}\left(\frac{\lambda_1}{\lambda_1 + \lambda}, \dots, \frac{\lambda_p}{\lambda_p + \lambda}\right) U^T \cdot y$

$Z = UD$

$=: S_\lambda$

$Z = X \cdot U^{-1}$

$\Rightarrow d(\text{ridge}) = d(S_\lambda) = d(U + D + \lambda I_p) = \frac{p}{\lambda + \dots}$

$$\Rightarrow \text{d}f(\text{ridge}_\lambda) = \text{tr}(S_\lambda) = \sum_{j=1}^p \frac{\lambda_j}{\lambda_j + \lambda} \quad (3.25)$$

$$\left(\begin{array}{l} \text{d}f(\text{ridge}_\lambda) < p, \text{ za } \lambda > 0, \text{ te} \\ \text{d}f(\text{ridge}_\lambda) \rightarrow 0, \text{ za } \lambda \rightarrow +\infty. \end{array} \right) \quad \text{d}f(LS)$$

[iako procenjujemo p parametara]

(Dz) Pokusite da je $\text{tr}\left(\frac{1}{\sigma^2} (\frac{1}{\sigma^2} + \lambda I)^{-1}\right) = \sum_{j=1}^p \frac{\lambda_j^2}{(\lambda_j + \lambda)^2}, \forall \lambda \geq 0.$
 (vidi Prop. 3.7).

3.3.2 Multikolinearnost

Uz model (3.5),

$\hat{\Sigma} = \frac{1}{n} X^T X$ = korelacijska matrica od $(x^{(1)}, \dots, x^{(p)}) \in \mathbb{R}^{n \times p}$

$$\text{Var}(\hat{\beta}_{LS}) = \frac{\sigma^2}{n} (\hat{\Sigma}^{-1})_{j,j}$$

\Rightarrow ako je $\hat{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$, $\text{Var}(\hat{\beta}_{LS}) = \frac{\sigma^2}{n} (\sigma_j^2)^{-1}$.

Općenito, vrijedi:

$$\text{Var}(\hat{\beta}_{LS}) = \frac{\sigma^2}{n} \text{Var}(x_j - \hat{x}_{j|I})^{-1} \quad (3.26)$$

gdje je $\hat{x}_{j|I} := \text{Proj}_{\text{span}(x_k)_{k \in I}}(x_j)$.

\Rightarrow ako $x_j \approx \sum_{k \in I} a_k x_k$ za neke $I \subseteq \{1, \dots, p\} \setminus \{j\}$ te $a_k \in \mathbb{R}$, $\text{Var}(\hat{\beta}_{LS})$ će biti veliko. "problem multikolinearnosti"

S druge strane, ako $\rho \ll n$, testna greška neće biti velika (ako je lin. model tačan, merenju).

40

mpr.) Prop. da $X_1 \approx X_2$ (dakle, $\rho(X_1, X_2) \approx 1$) te ujedn.:

$$Y = \underbrace{X_1 + X_2}_{\approx cX_1 + (2-c)X_2} + \varepsilon \quad (\text{dakle, } \beta_0 = (1, 1))$$

Intuitivno

\Rightarrow $RSS(c, 2-c) \approx RSS(1, 1)$, $\forall c \in \mathbb{R}$, a argmin $RSS(\beta)$ $\beta \in \mathbb{R}^2$ jeko ovise o X i y , tj. T !

\hookrightarrow ridge regresija penalizira više veličine $(c, 2-c)$ za $|c| \gg 0$, pa $\text{Var}(\hat{\beta}_j^r) < \text{Var}(\hat{\beta}_j^{ols})$.

\rightarrow ridge - kolinearnost. R

Kako prepoznati multikolinearnost

(1) pogledati korelacije među korelativama

(2) Variance Inflation Factor:

$$VIF(\beta_j) := \frac{(\hat{\Sigma}^{-1})_{jj}}{\lambda_j^{-1}} \quad (3.27)$$

$$= \frac{\text{varijanca od } \hat{\beta}_j^{ols}}{\| \text{u slučaju } \forall_j^T X_n = 0, \forall n \neq j$$

\hookrightarrow ofudeno je li $VIF(\beta_j) \geq 5$ ili 10. (hećenodika)

(3) Ako $\lambda_k = \text{suvojstr. vrijedn. od } \hat{\Sigma} \approx 0$, [k-ta po veličini]

$$\Rightarrow 0 \approx \lambda_k = \widehat{\text{Var}}\left(\frac{\hat{\beta}_k}{\sqrt{\lambda_k}}\right) = \widehat{\text{Var}}\left(\sum_{j=1}^p X_j \cdot V_{jk}\right)$$

$$\Rightarrow \sum_{j=1}^p X_j V_{jk} \approx \text{const!}$$

PA