

3) Regularizacija linearnog modela [ridge regresija]

Nastajanje

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

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

↙ 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 (ujednost. koeficijenti 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_{n0}, \dots, x_{n0})

Lin. model: biramo $\hat{f} = d(w)$ 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) kvadratne koeficijente.]

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

$$\hat{\beta}^{ls} = (\hat{\beta}_0^{ls}, \dots, \hat{\beta}_p^{ls}) := \underset{\beta \in \mathbb{R}^{p+1}}{\operatorname{argmin}} \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - f_\beta(x^{(i)}))^2}_{=: \operatorname{RSS}(\beta)} \quad (3.2)$$

$$= \underset{\beta}{\operatorname{argmin}} \|y - X\beta\|_2^2$$

Pretpostavimo da su $x_0, \dots, x_p \in \mathbb{R}^m$ lin. nezavisni, tj. da je 24

$X^T X \in \mathbb{R}^{(p+1) \times (p+1)}$ — poz. definitna (simetrična) matrica.
(misao je m ≥ 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 $\hat{\Sigma} := \frac{1}{n} X^T X$.

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

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

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

Nap.1 U slučaju da x_0, \dots, x_p nisu lin. nezavisni (npr. $m < p+1$ ili $x_j = 0 \cdot x_i$) $\hat{\beta}^{OLS}$ nije jedinstven. Ipak, \hat{y} je i dalje ortogonalna projekcija od y na span(x_0, \dots, x_p) ⊥

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

[masovni pretpostavni model za (X, Y)]

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

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

pri čemu su [nodi jedinstvenosti]

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

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

$$Y := \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = X \beta_* + \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix} =: X \beta_* + \varepsilon, \quad (3.6)$$

[kiks]
↓

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

⇒ (3.3) $\hat{\beta}^{OLS} = \beta_* + \left[\frac{1}{n} \Sigma^{-1} X^T \varepsilon \right]$ jedini slučajni

• $E[\hat{\beta}^{OLS}] = \beta_*$

• $\text{Var}(\hat{\beta}^{OLS}) = \frac{\sigma^2}{n} \Sigma^{-1}$

(3.7)

Tipično

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

3.2.1 Testna greška

Neka je $\tilde{Y} := X \beta_* + \tilde{\varepsilon}$, pri čemu su $\varepsilon_i \sim \tilde{\varepsilon}$ (tj. $y_i \sim \tilde{y}_i$) nezavisni, testni skup.

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

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

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

Prop. 3.11 (BVT) Ako je $\hat{\beta} = \hat{\beta}(Y, X)$ procjenitelj za β_0 , 26

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

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

Dokaz

(1) Iz (3.6) imamo, $\forall \beta \in \mathbb{R}^{p+1}$, [σ^2 + reducirana greška]

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

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

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

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

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

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

□

Prop. 3.2

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

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

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

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

↑
nije slučajno

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

□

Alternativ,

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

$$= \hat{\beta}_{OLS} - \beta_0 = \text{tr}(\epsilon^T H \epsilon)$$

$$E[\epsilon \epsilon^T] = \sigma^2 \mathbb{I}_n$$

↓
+ tr(H) = p+1

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

(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.

Nadalje, iz Prop. 2.5 imamo $(E[L_T(\hat{\beta})] = E_Y[L^*(\hat{\beta})])$,

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

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

[Specijalno, za $\hat{\sigma}^2$ iz (3.8) imamo $E[\hat{\sigma}^2] = \sigma^2$.]

(ii) Dakle, čak i kad je lin. model tačan (pa je pristranost = 0), varijanca (a sumirane i ukupne greške) je velika (u odnosu na σ^2) ako je $\frac{p+1}{n}$ veliko!

Uostalom, svaka dodatna kovarijata x_j povećava grešku
za $\frac{\sigma^2}{n}$ bez obzira je li $\beta_j \approx 0$ ili
 $x_j \approx x_i$ za $i \neq j$.



3.3 Ridge regresija

Pr | (Motivacija)

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

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

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

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

\Leftrightarrow

$n \geq \frac{3\sigma^2}{\theta^2}$

velike obr $\theta \approx 0$.

"L2 regularizacija"

Za $\lambda \geq 0$ defin. $\hat{\beta}_\lambda$

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

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

(3.13)

me uključuje β_0

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

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

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

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

centriranje

Zbog Prop. 3.4 u nastavku prepostavljamo da je

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

te rješavamo

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

$\nwarrow \mathbb{R}^{n \times p}$

(vidi (3.24))

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

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

"nepošteni" akteri koeficijenti nisu u istim mjerama jedinica.

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. 1 (i) za $\lambda > 0$, $\hat{\beta}^r$ je uvijek dobro definiran (dobre, 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) koefic. prema 0.

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

$$E \left[L^{in}(\hat{\beta}_{\lambda}^r) \right] = d^2 + \overbrace{\lambda^2 \beta_a^T (\hat{\Sigma} + \lambda I_p)^{-2} \hat{\Sigma} \beta_a}^{\text{bias}} + \underbrace{\frac{d^2}{n} \operatorname{tr} \left[\hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2} \right]}_{\text{varijanca}}. \quad (3.17)$$

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

$$Y = X\beta_0 + \varepsilon$$

iz Prop. 3.1, $\mathbb{E}[\| \hat{\beta}_\lambda^r - \beta_0 \|^2] = B + V$, gdje

$$B = \mathbb{E}[\| \hat{\beta}_\lambda^r - \beta_0 \|^2]$$

$$\begin{aligned} \mathbb{E}[\hat{\beta}_\lambda^r] &= \frac{1}{n} (\hat{\Sigma} + \lambda I_p)^{-1} X^T X \beta_0 + \frac{1}{n} (\hat{\Sigma} + \lambda I_p)^{-1} X^T \mathbb{E}[\varepsilon] \\ &= \frac{1}{n} (\hat{\Sigma} + \lambda I_p)^{-1} X^T X \beta_0 + \frac{1}{n} (\hat{\Sigma} + \lambda I_p)^{-1} X^T \mathbb{E}[\varepsilon] \\ &= \beta_0 - \lambda (\hat{\Sigma} + \lambda I_p)^{-1} \beta_0 \end{aligned}$$

$$\begin{aligned} \Rightarrow B &= \mathbb{E}[\| \lambda (\hat{\Sigma} + \lambda I_p)^{-1} \beta_0 \|^2] \\ &= \lambda^2 \beta_0^T (\hat{\Sigma} + \lambda I_p)^{-1} \frac{1}{n} (\hat{\Sigma} + \lambda I_p)^{-1} \beta_0 \\ &= \lambda^2 \beta_0^T (\hat{\Sigma} + \lambda I_p)^{-2} \frac{1}{n} \beta_0 \end{aligned}$$

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

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

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

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

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

now done

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

Uput-1 (i) Pokazati 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. vrednosti od $\hat{\Sigma}$.

\Rightarrow
 $\forall \lambda > 0$

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

(ii)

pristrenost = $\beta_\alpha^\top \left(\frac{\hat{\Sigma}}{\lambda} + I_p \right)^{-2} \hat{\Sigma} \beta_\alpha \rightarrow \beta_\alpha^\top \hat{\Sigma} \beta_\alpha = \|\beta_\alpha\|_{\hat{\Sigma}}^2$,
($\neq 0!$)

kada $\lambda \rightarrow \infty$. \square

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

$$E \left[L^{in}(\hat{\alpha}_{\beta_\alpha}^{\lambda_\alpha}) \right] \geq E \left[L^{in}(\hat{\alpha}_{\beta_\alpha}^{\lambda_\alpha}) \right]$$

\hookrightarrow u prekos, a odabiremo CV metodom.

[λ_α ne možemo eksplicitno označiti jer onisi su reprezentativni relacijama.]

3.3.1 Kako vidje smanjuje koeficijente?

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

$$X = U \cdot D \cdot V^T, \quad (3.18)$$

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

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 ortogonalizirani.

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

Zbog pref. $\bar{x}_j = 0, \forall j=1, \dots, p,$

$$\hat{\Sigma} = \frac{1}{n} X^T X = \text{matrica kovarijancijska matrica}$$

$\in \mathbb{R}^{p \times p}$ od $x^{(1)}, \dots, x^{(n)}$.

$\hat{\Sigma} \stackrel{(3.18)}{=} V \frac{D^2}{n} V^T,$

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

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

tzv. svojstveni glaovni komponenti (od X)

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

Nadalje, vektore

$$Z_j := X v_j = \begin{bmatrix} v_j^T x^{(1)} \\ \vdots \\ v_j^T x^{(n)} \end{bmatrix} \in \mathbb{R}^{n \times 1}, j=1, \dots, P \quad (3.145)$$

zove se 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) Vnjeti, $\forall j=1, \dots, P$,

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

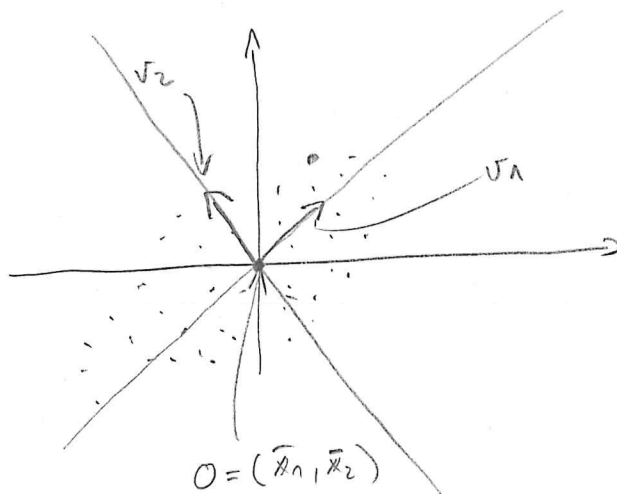
↑
varijanca varijena od $v^T x^{(i)}, i=1, \dots, n$

uz uvjet

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

[bez dokaza → Primjenjena statistika]

prim. $P=2$



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

Neka je

$$Z := X V = \begin{bmatrix} | & & | \\ z_{11} & \dots & z_{1p} \\ | & & | \end{bmatrix} \in \mathbb{R}^{n \times p}, \quad (3.20)$$

zapravo podataka u bazi v_1, \dots, v_p .

(3.18) $\Rightarrow Z = U \cdot D, \quad t_j$

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

(3.21)

Specijalno,

- $\|z_j\|_2^2 = z_j^T z_j = d_j^2, \quad \forall j,$

(3.22)

te

$$\widehat{\text{Var}}(z_j) = \frac{d_j^2}{n} = \lambda_j, \quad \forall j.$$

$$[z_j = (X^T)_j \cdot v_j = 0, \quad \forall j=1, \dots, p]$$

- $z_j^T z_k = 0 = \widehat{\text{Cov}}(z_j, z_k), \quad \forall j \neq k$

(3.23)

$$\Rightarrow \text{span}(x_1, \dots, x_p) = \text{span}(z_1, \dots, z_p).$$

→ Ridge regresija i glavne komponente

Pretp. da je $\widehat{\text{Cov}}(x_j, x_k) = 0, \quad \forall j \neq k$

$$\Rightarrow X^T X = \text{diag}(\|x_1\|_2^2, \dots, \|x_p\|_2^2)$$

Isto kao kod regresije od y samo na x_j !

te

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

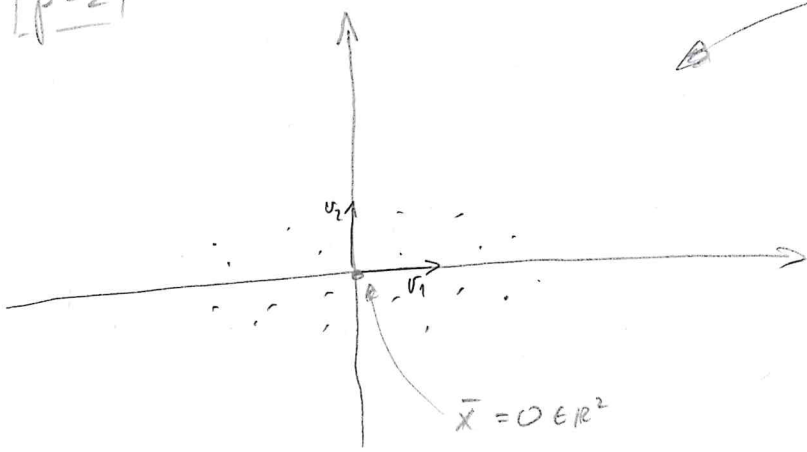
veće smanjenje ako je $\|x_j\|_2^2 = n \cdot \widehat{\text{Var}}(z_j)$ manja!

(3.24)

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

(3.16) $\langle 1, \text{ za } \lambda > 0 \rangle!$

mp.1 | $p=2$



$\text{Var}(x_1) > \text{Var}(x_2)$
te u više velik

Intuitivno, ako prilagodavamo lin. pljku (y-os je tražena koordinata), $\text{Var}(\hat{\beta}^n)$ (a samim time i varijanca u $E[L^{\text{in}}(\hat{\beta}^n)]$) bit će velika jer prognoza manjka u smjeru v_2 jako ovise o T.

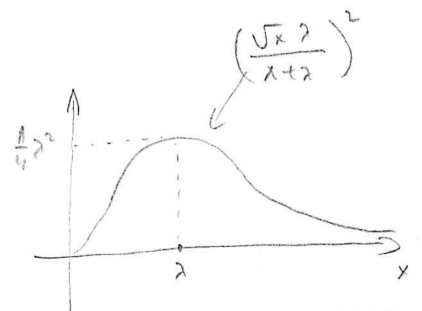
Ridge regresija smanjuje tu varijancu t.d. smanjuje prognozu manjka u smjerovima v_1 i v_2 , i to više u smjeru (v_2) .

→ pristaoštvu ovise o λ te o tome je li stvarni manjka u smjeru v_2 (t.j. β_2) velik ili ne!

↳ Iz $\hat{\Sigma} = \text{diag}(\hat{\text{Var}}(x_j); j=1, \dots, n)$ imamo da je pristaoštvu u $E[L^{\text{in}}(\hat{\beta}^n)]$ jednaka (Dt)

$$\sum_{j=1}^p \left(\frac{\sqrt{\Delta_j} \cdot \lambda}{\Delta_j + \lambda} \right)^2 (\beta_j^*)^2$$

max za $\Delta_j = \lambda$.
 $\Delta_j \rightarrow +\infty$
 $\Delta_j \rightarrow 0$



→ ridge - simulacijski - primjer. pdf

Općenito, $V\beta \in \mathbb{R}^p$

$$X = UDV^T = ZV^T$$

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

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



$$Z^T Z = \text{diag}(d_1^2, \dots, d_p^2), \quad \text{gdje}$$

(3.22), (3.23)

$$\Rightarrow \hat{d}_j^r = \frac{d_j^2}{d_j^2 + n\lambda} \cdot \hat{d}_j^{\text{ls}}, \quad j = 1, \dots, p$$

(3.24)

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

Općeniti, vidjele smo prije da su projekcije, ali nakon projekcije podataka na v_1, \dots, v_p , te je smanjenje varijance za \hat{d}_j^{ls} za koje je $\widehat{\text{Var}}(\hat{z}_j) = \frac{d_j^2}{n} = \lambda_j$ manje.

Prop. 3.21 $\text{diag}(vidje) = ?$

Imamo $\forall \lambda \geq 0$, (3.16) $\text{diag}(d_1^2, \dots, d_p^2)$

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

$$= U \cdot \text{diag}\left(\frac{d_1^2}{d_1^2 + n\lambda}, \dots, \frac{d_p^2}{d_p^2 + n\lambda}\right) U^T \cdot y$$

$Z = UD$

$$=: S_\lambda$$

$Z \cdot \beta = UD \cdot \beta$

$= U \cdot D \cdot \beta = U \cdot \beta$

$(U \cdot \beta)^T = \beta^T U^T$

$(U \cdot \beta)^T \cdot y = \beta^T U^T y$

$$\Rightarrow \text{d}t(\text{ridge}_\lambda) = \text{tr}(S_\lambda) = \text{tr}(\underbrace{U^T U}_{=I_p} \text{diag}(\cdot)) \quad [38]$$

$$= \sum_{j=1}^p \frac{d_j^2}{d_j^2 + n\lambda} \quad \left(= \sum_{j=1}^p \frac{\lambda_j}{\lambda_j + \lambda} \right) \quad (3.25)$$

$\lambda_j = \frac{d_j^2}{n}$

$\text{d}t(\text{ridge}_\lambda) < p$, za $\lambda > 0$, te
 $\text{d}t(\text{ridge}_\lambda) \rightarrow 0$, za $\lambda \rightarrow +\infty$.

[iako procenjujemo p parametara]

(D2) Pokazite da je $\text{tr}(\hat{\Sigma}^{-2} (\hat{\Sigma} + \lambda I)^{-2}) = \sum_{j=1}^p \frac{\lambda_j^2}{(\lambda_j + \lambda)^2}$, $\forall \lambda \geq 0$.
 (vidi Prop. 3.7).

3.3.2 Multikolinearnost

(K2 model (3.5),

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

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

\Rightarrow ako je $\hat{\Sigma} = \text{diag}(\underbrace{\sigma_1^2}_{\text{Var}(x_1)}, \dots, \sigma_p^2)$, $\text{Var}(\hat{\beta}_{OLS}) = \frac{\sigma^2}{n} (\sigma_j^2)^{-1}$.

Općenito, vrijedi:

[varijance]

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

gdje je $\hat{x}_{j-1-j} := \text{Proj}_{\text{span}(x_k, k \neq j)}(x_j)$.

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

S druge strane, ako $\rho \ll n$, testna greška neće biti velika (ako je lin. model tačan, mereno). (30)

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

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

Intuitivno

\Rightarrow $RSS(c, 2-c) \approx RSS(1, 1)$, $\forall c \in \mathbb{R}$, a argmin $RSS(\beta)$ $\beta \in \mathbb{R}^2$ jeku ovim X i y !

\hookrightarrow ridge regresija penalizira više rektore $(c, 2-c)$ za $|c| \rightarrow 0$,

$$\text{pa } \text{Var}(\hat{\beta}_j^r) \stackrel{(\Leftarrow)}{<} \text{Var}(\hat{\beta}_j^{OLS}).$$

\rightarrow ridge - kolinearnost. R

Kako prepoznati multikolinearnost

(1) pogledati korelacije među korekpondentima

(20) Variance Inflation Factor:

$$VIF(\beta_j) := \frac{(\hat{\Sigma}^{-1})_{jj}}{s_j^{-2}} \quad (3.27)$$

$$= \frac{\text{varijanca od } \hat{\beta}_j^{OLS}}{\| \cdot \| \text{ u slučaju } \sum_{j \neq k} x_k = 0, \forall k \neq j}$$

\hookrightarrow sledeno je li $VIF(\beta_j) \geq 5$ ili 10. (heuristika)

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

$$\Rightarrow 0 \approx \lambda_k = \widehat{\text{Var}}\left(\frac{z_k}{n}\right) = \widehat{\text{Var}}\left(\sum_{j=1}^p x_j \cdot V_{jk}\right)$$

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