

② Procjena greške i odabir modela
 [Model selection]

Za dani T , zanima nas procjena greške

$$L(\hat{f}_T) = \mathbb{E}[L(Y, \hat{f}_T(x))] \quad (2.1)$$

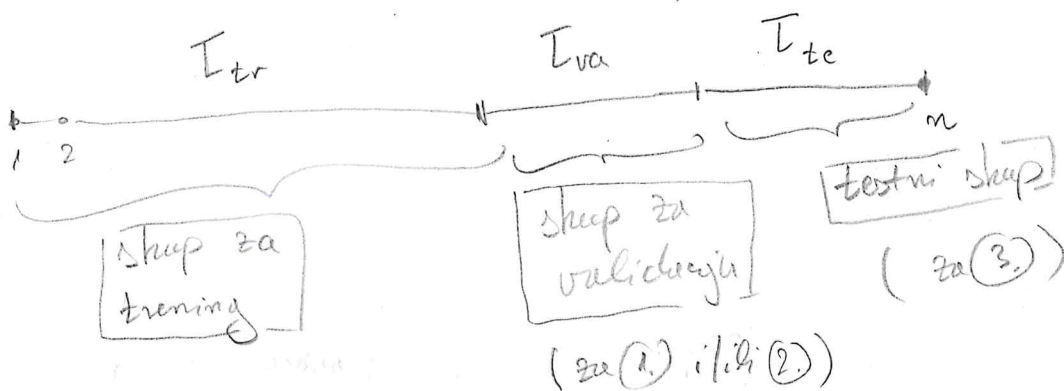
↳ hiperparametar koji kontrolira kompleksnost modela
 (npr. $\alpha = \frac{1}{n}$ ili h u boost metod.)

Ciljevi:

- ① Odabir hiperparametra α ;
- ② Odabir između različitih modela; [Model selection]
- ③ Procjena testne greške odabranog modela. [Model assesment]

Ako je n velik, najbolji pristup je (na slučajnoj način)

podijeliti T :



1. T_{tr} → prilagodba raznih modela i/ili istog modela za različite α [?]
2. T_{va} → za svaki model ili α procjenjujemo testnu grešku (2.1)
 te npr. izaberemo onaj s najmanjom procjenjenom greškom

3. $T_{te} \rightarrow$ prognozujeemo testnu gresku (2.1) odobrenog modela.

Def. 2.11 Zashto trebamo T_{te} ?

npr. ako je T_{ra} realizacija od njez uvrstva

$$T_{ra} = \{(\tilde{X}^{(i)}, \tilde{Y}_i) : i=1, \dots, m\}$$

za (X, Y) , za nekak funkciju λ -ju $\lambda: \mathbb{R}^p \rightarrow \mathbb{R}$,

$$\underline{\underline{L_{Tra}(\lambda) = \frac{1}{m} \sum_{i=1}^m L(\tilde{Y}_i, \lambda(\tilde{X}^{(i)}))}} \quad (2.2)$$

$\stackrel{d}{\approx} N(\underbrace{E[L(Y, \lambda(X))]}_{= L(\lambda)}, \frac{\text{Var}(L(Y, \lambda(X)))}{m})$
CGT

[Stocajna varijabla]
 $\lambda \lambda!$

ako je m dovoljno velik.

pa, ako nećemo $L_{Tra}(\hat{\lambda}_\lambda)$ za sve λ te

$$\hat{\lambda} = \underset{\lambda}{\text{argmin}} L_{Tra}(\hat{\lambda}_\lambda), \quad (2.3)$$

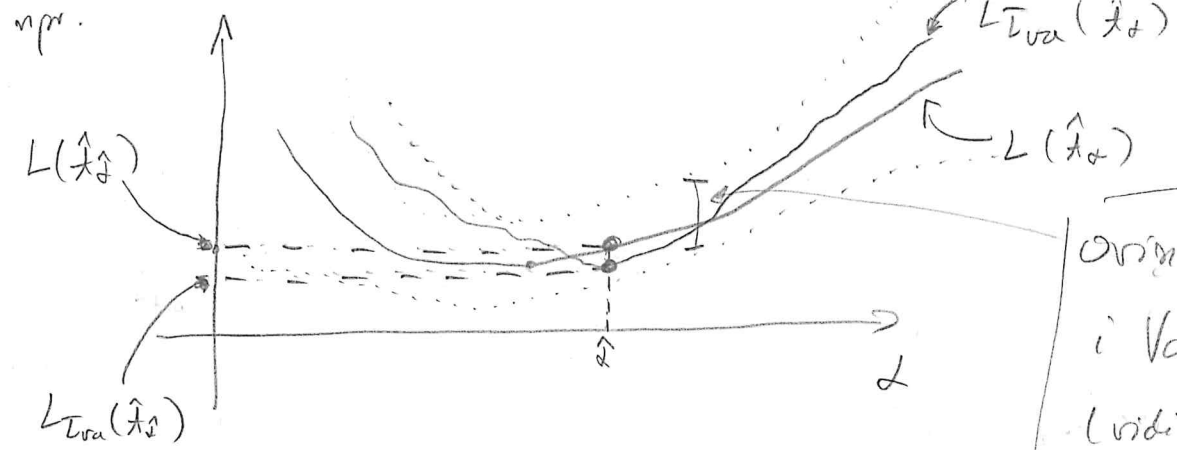
imamo

$$E[L_{Tra}(\hat{\lambda}_\lambda)] \stackrel{\ominus}{\leq} \min_{\lambda} E[L_{Tra}(\hat{\lambda}_\lambda)]$$

$\leq L_{Tra}(\hat{\lambda}_\lambda), \forall \lambda$
 po (2.3)

$$\stackrel{\ominus}{\leq} L(\hat{\lambda}_\lambda) \quad \leftarrow \text{[Stocajna varijabla]}$$

$$E[L(\hat{\lambda}_\lambda) - L_{Tra}(\hat{\lambda}_\lambda)] \geq 0! \quad \text{[pristranost]}$$



Ornisi σ (m)
 $i \text{ Var}(L(Y, \hat{x}_i(x)))$
 (vidi (2.2))

[veći m , manja greška]

Slučaj, ako imamo više različitih modela $\rightarrow \hat{x}^{(1)}, \dots, \hat{x}^{(k)}$ te

$$\hat{i} := \underset{i=1, \dots, k}{\text{argmin}} L_{Tra}(\hat{x}^{(i)}),$$

imamo

$$E[L_{Tra}(\hat{x}^{(\hat{i})})] \leq L(\hat{x}^{(\hat{i})}).$$

[Malo je komplikovanije u slučaju kada biramo $\hat{x}^{(i)}$ i \hat{i} .]

2.1 Umaksna validacija (CV)

Kada (m) nije velik, tj. T_{tr}, T_{va} i T_{te} su premaleni, "izbacujemo" T_{va} te odaberemo modela i/ili hiperparametara na osnovu T_{tr} .

↳ u tom slučaju CV predstavlja najjednostavniju i najpopularniju metodu.

"k-fold CV": (neka je $\mathcal{T} := \mathcal{T}_{tr}$)

- za dani $k \in \{2, 3, \dots, n\}$ (na sledećem načinu) podjelimo \mathcal{T} na k približno jednakih disjunktih "dijelova" ("folds")

$\mathcal{T}_1, \dots, \mathcal{T}_k$

→ neka je $n_j := |\mathcal{T}_j| = \# \text{ elemenata u } \mathcal{T}_j, j=1, \dots, k$

→ pretp. $n_k := \frac{n}{k} \in \mathbb{N} \Rightarrow \boxed{n_j = n_k}, \forall j=1, \dots, k$
veličina svakog bloka

- $\forall j=1, \dots, k$, računamo [onih's algoritmu]

$$CV_j^{(h)} = CV_j^{(h)}(\hat{A}) := \frac{1}{n_k} \sum_{(x^{(i)}, y_i) \in \mathcal{T}_j} L(y_i, \hat{A}^{-j}(x^{(i)})), \quad (2.4)$$

pri čemu je

$$\hat{A}^{-j} := \hat{A}(\mathcal{T}^{-j})$$

za

$$\mathcal{T}^{-j} := \bigcup_{i \neq j} \mathcal{T}_i$$

ovim \mathcal{T}_j

- k-CV procjena testne greške od $\hat{A} = \hat{A}(\mathcal{T})$ je

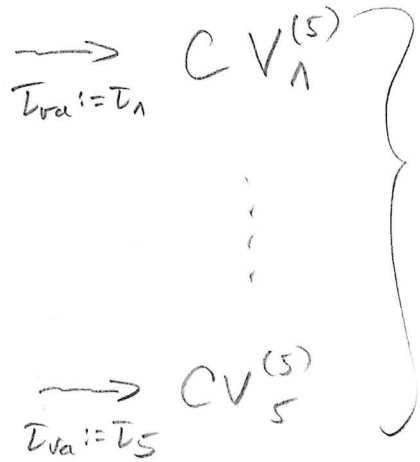
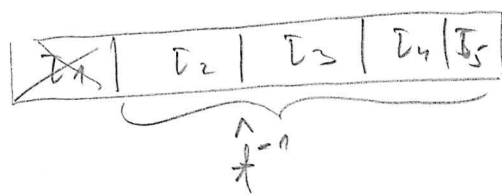
$$CV^{(h)} = L_{cv}^{(h)}(\hat{A}) := \frac{1}{k} \sum_{j=1}^k CV_j^{(h)} \quad (2.5)$$

$$= \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{A}^{-j(i)}(x^{(i)})), \quad (2.6)$$

pri čemu je $j(i) \in \{1, \dots, k\}$ t.d. $(x^{(i)}, y_i) \in \mathcal{T}_j$.

[(2.6) je dobro def. i u slučaju kada $\frac{n}{k} \notin \mathbb{N}$, tj. blokovi su različite dužine.]

mpm. $h=5$



$L_{CV}^{(h)}(\hat{\beta})$

↳ u j -tom koraku β_j igra ulogu shepa za validaciju (tj. nezavisnog testnog shepa).

Nap. 2.21 Ispostavlja se da $L_{CV}^{(h)}(\hat{\beta})$ bolje procjenjuje $E_T[L(\hat{\beta})]$

nego $L(\hat{\beta})$ (za danu $T=\bar{T}$) (vidi ESL, pogledajte 7).

Ipak, ako $L_{CV}^{(h)}$ koristimo za odabir modela / hiperparametara tipično izabiremo model / hiperparameter koji ima manju testnu grešku $L(\hat{\beta})$. [CV je još uvijek tečtški nedovoljno istražen!] □

Kako odabrati h ?

↳ $|\bar{T}-s| = n - r_h = \frac{k-1-h}{k} \cdot n < n$

⇒ tipično imamo $E_T[L_{CV}^{(h)}(\hat{\beta})] \geq E_T[L(\hat{\beta})]$

za male h (odnos velike r_h)

pristranst!

S druge strane, u slučaju $k=m$, tj. $T^{\rightarrow} = T \setminus \{(x^{(j)}, y_j)\}, j=1, \dots, n$, imamo manju pristupaost, ali budući da su $T^{-j}, j=1, \dots, n$ vrlo bliski, procjene $CV_j^{(n)}, j=1, \dots, n$ su pozitivne korelacije pa je $Var_T(L_{CV}^{(h)}(\hat{\alpha}))$ veća.

[opet BVT!]

→ tzv. "leave-one-out CV" (LOOCV)

[→ za velike n , $L_{CV}^{(n)}$ može biti zahtjevan za izračunati.]

U praksi najčešće $k=5$ ili 10 te $k=n$.

[Također ne postoji prava teorijska rezultata.]

Nap 2.3 Umjesto $\hat{\alpha} = \hat{\alpha}_{min}$ za $\hat{\alpha}_{min} := \arg\min_{\alpha} L_{CV}^{(h)}(\hat{\alpha})$ često se uzima $\hat{\alpha}$ koji odgovara najjednostavnijem modelu t.d.

$$L_{CV}^{(k)}(\hat{\alpha}) \leq L_{CV}^{(h)}(\hat{\alpha}_{min}) + \widehat{SE}(\hat{\alpha}_{min}),$$

gdje je

$$\widehat{SE}(\hat{\alpha}) := \sqrt{\frac{Var(CV_1^{(h)}, \dots, CV_n^{(h)})}{k}} \quad \leftarrow \text{varijanca varijanci}$$

(2.7)

(naivni) procjenitelj za standardnu grešku $\sqrt{Var_T(L_{CV}^{(h)}(\hat{\alpha}_*))}$ procjenitelja $L_{CV}^{(h)}(\hat{\alpha}_*)$.

• (2.7) ima smisla kada je k velik i $CV_1^{(h)}, \dots, CV_n^{(h)}$ približno nekorelacijski! (2.7) podcijenjuje stvarnu std. grešku)

• tzv. "one SE rule"

• analogno pravilo možemo primijeniti kod otlabine modela

• Ozračar princip: između jednakih preciznih modela

(do na jednu SE), biramo majjednostavniji.

→ CV- ilustracija. R

Pr. 2.41 Pretp. da je $p = 5000$, a $m = 100$ (dakle, $p \gg m$),

te $Y \in \mathbb{R}^m$. Naš algoritam je

1. Na temelju T (izračunaj uzorke korelacije

$$P((x_{1j}, y_1), (x_{2j}, y_2), \dots, (x_{mj}, y_m)),$$

j -te kovarijate i odabira, $j = 1, \dots, p$.

→ Ostavi samo $m = 100$ kovarijata sa majjednom korelacijom

2. Koristeći samo tih m kovarijata i neku metodu generiraj $\hat{\lambda}$.

→ Kako izabrati CV algoritam za procjenu greske?

→ i 1. mora biti uključen u srži korek CV algoritma!

(vidi ESL, poglavlje 7.10.2)

Npr., neka je

• X_1, \dots, X_p mjd $\sim N(0,1)$ nezavisno

• $Y \sim \begin{pmatrix} 0 & 1 \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix}$ nezavisno od X_1, \dots, X_p

⇒ $\forall \lambda: \mathbb{R}^p \rightarrow \mathbb{R}^m$ vrijedi

$$L(\hat{t}) = P(Y \neq \hat{t}(X))$$

prop. 0-1 gubitak

$$= \frac{1}{2} P(\hat{t}(X)=1 | Y=0) + \frac{1}{2} P(\hat{t}(X)=0 | Y=1)$$

$$= \frac{1}{2} P(\hat{t}(X)=1) + \frac{1}{2} P(\hat{t}(X)=0) = \left[\frac{1}{2} \right]$$

mer. od Y i X

[manjima!]

\Rightarrow za bilo koju metodu $i \in \mathcal{T}$,

$$L(\hat{t}) = \frac{1}{2} = E_{\mathcal{T}}[L(\hat{t})].$$

\rightarrow CV-primer. R

Nap. ukoliko radimo korelaciju medu nekim metodom neregularnog učenja (npr. analiza glavnih komponenti), tj. ne koristimo y_i -ove, taj korek nije potrebno uvesti u CV algoritom!

2.2) Stepnjeni slobode (degrees of freedom) \rightarrow (dt)

\hookrightarrow želim upoređivati kompleksnost različitih metoda

Pretp. da je model kao u (1.10), $Y = \hat{t}(X) + \epsilon$, te nadalje da je

$$Y_i = \underbrace{\hat{t}(x^{(i)})}_{=: \mu_i} + \epsilon_i, \quad i=1, \dots, n$$

uz

• $\epsilon_1, \dots, \epsilon_n$ n. nezavisne t.d. $E[\epsilon_i] = 0, \text{Cov}(\epsilon_i, \epsilon_j) = 0, \forall i \neq j$

• $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^p$ fiksni $\hookrightarrow T = \{(x^{(i)}, Y_i) : i=1, \dots, n\}$
("fixed-design")

(alternativno, mogli smo u nastavku uvijek brati $x^{(i)} = x^{(i)}$, $i=1, \dots, n$)

uz $Y = (Y_1, \dots, Y_n)^T \in \mathbb{R}^{n \times 1}$, $\mu = (\mu_1, \dots, \mu_n)^T \in \mathbb{R}^{n \times 1}$ te

$\mathcal{E} = (\epsilon_1, \dots, \epsilon_n)^T \in \mathbb{R}^{n \times 1}$, (2.8) postaje

$Y = \mu + \mathcal{E}$, $E[\mathcal{E}] = 0$, $\text{Cov}(\mathcal{E}) = \sigma^2 I_n$

↑
kovarijancijska
matrica

Za metodu $A: T \rightarrow \hat{A}(T)$ definiramo

$$d_A(A) := \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(\hat{Y}_i, Y_i)$$
 (2.9)

jedina slučajnost dolazi od $\epsilon_1, \dots, \epsilon_n$

pri čemu je $\hat{Y}_i := \hat{A}(x^{(i)})$, $i=1, \dots, n$.

[dobre, ako se "više" prilagođavamo podacima, Cov pa osemim time i d_A raste!]

Motivacija za (2.9): usporediti čemo

• $E \left[\underbrace{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}_{L_T(\hat{A})} \right]$ ← očekivana greška na skupu za treniranje

• $E \left[\underbrace{\frac{1}{n} \sum_{i=1}^n (\tilde{Y}_i - \hat{Y}_i)^2}_{L_{\tilde{T}}(\hat{A})} \right]$ ← tzv. (očekivana) "in-sample" testna greška

gdje $\tilde{T}_i = \{ (x^{(i)}, \tilde{Y}_i) : i=1, \dots, m \}$ testni skup nezavisan

od $T (\tilde{Y}_i = \mu_i + \tilde{\epsilon}_i, i=1, \dots, m), (\tilde{\epsilon}_i)_{i=1, \dots, m}$ nez. od $(\epsilon_i)_{i=1, \dots, m}$.

Prp. 2.5 | Vrijedi:

$$E[L_T(\hat{\beta})] - E[L_T(\hat{\beta}^*)] = \frac{2d^2}{n} \cdot d\beta(A) \quad (2.10)$$

ter. očekivani "optimalizam"
od $L_T(\hat{\beta}^*)$

Dokaz:

Rijeva strana u (2.10) = $E \left[\frac{1}{n} \sum_{i=1}^m \left\{ \tilde{Y}_i^2 - 2\tilde{Y}_i \hat{Y}_i + \hat{Y}_i^2 - Y_i^2 + 2Y_i \hat{Y}_i - \hat{Y}_i^2 \right\} \right]$

$$= \frac{1}{n} \sum_{i=1}^m \left\{ \underbrace{E[\tilde{Y}_i^2]}_{= \text{Var}(\tilde{Y}_i) + \mu_i^2 = d^2 + \mu_i^2} - 2 \underbrace{E[\tilde{Y}_i \cdot \hat{Y}_i]}_{= E[\tilde{Y}_i] E[\hat{Y}_i] = \mu_i} - \underbrace{E[\hat{Y}_i^2]}_{= \text{Var}(\hat{Y}_i) + \mu_i^2 = d^2 + \mu_i^2} + 2 E[Y_i \hat{Y}_i] \right\}$$

$$= \frac{2}{n} \sum_{i=1}^m \left\{ E[Y_i \hat{Y}_i] - \mu_i \cdot E[\hat{Y}_i] \right\} = \frac{2}{n} \sum_{i=1}^m \text{Cov}(Y_i, \hat{Y}_i) = \frac{1}{n} \text{Var}(Y_i) = \frac{1}{n} d^2$$

Pr. 2.6 |

(i) Ako je $\hat{Y}_i := \bar{Y} = \frac{1}{n} \sum_{i=1}^m Y_i, \forall i \Rightarrow d\beta(A) = \frac{1}{d^2} \sum_{i=1}^m \text{Cov}(Y_i, \bar{Y}) = \boxed{1}$

(ii) Ako je $\hat{Y}_i := Y_i, \forall i \Rightarrow d\beta(A) = \frac{1}{d^2} \sum_{i=1}^m \text{Var}(Y_i) = \boxed{m}$

(iii) Kod lin. regresije metodom najmanjih kvadrata, $d\beta(A) = P$ (kasnije)

(iv) Kod kNU regresije, $d\beta(A) = \frac{n}{k}$ (Dž)

Svi primjeri iz Pr. 2.5 su oblika

$$\hat{Y}_i = \sum_{j=1}^n S_{ij} Y_j, \quad i=1, \dots, n \quad (2.11)$$

gdje je $S = (S_{ij} : i, j \in \{1, \dots, n\}) \in \mathbb{R}^{n \times n}$ (deterministička)

matrica čiji elementi ovise samo o $x^{(1)}, \dots, x^{(n)}$.

Prop. 2.71 Ako je A t.d. ujedr. (2.11),

$$dA(A) = \text{tr}(S) \quad (2.12)$$

Dokaz

$$\text{Cov}(Y_i, \hat{Y}_i) = \text{Cov}\left(Y_i, \sum_{j=1}^n S_{ij} Y_j\right)$$

$$= \sum_{j=1}^n S_{ij} \underbrace{\text{Cov}(Y_i, Y_j)}_{=0 \text{ } i \neq j} = S_{ii} \cdot \underbrace{\text{Var}(Y_i)}_{=\sigma^2}, \quad \forall i=1, \dots, n$$

S_{ij} nisu slučajni

Prop. 2.71 Iz (2.11) slijedi da je statistika

$$T = L_T(\hat{A}) + \frac{2\sigma^2}{n} dA(A) \quad (2.13)$$

nepristran procjenitelj za $\mathbb{E}[L_T(\hat{A})]$

$\rightarrow T$ možemo koristiti za odabir modela.

\rightarrow Kod lin. regresije metodom najmanjih kvadrata,

$$(2.14) \quad T = L_T(\hat{A}) + \frac{2\sigma^2}{n} \textcircled{P} \quad \leftarrow \text{tzv. CP statistika} \quad (2.14)$$

(koristi se za odabir varijabli)