

② Procjena greške i odabir modela  
 [Model selection]

Za dani  $T$ , zanima nas procjena greške

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

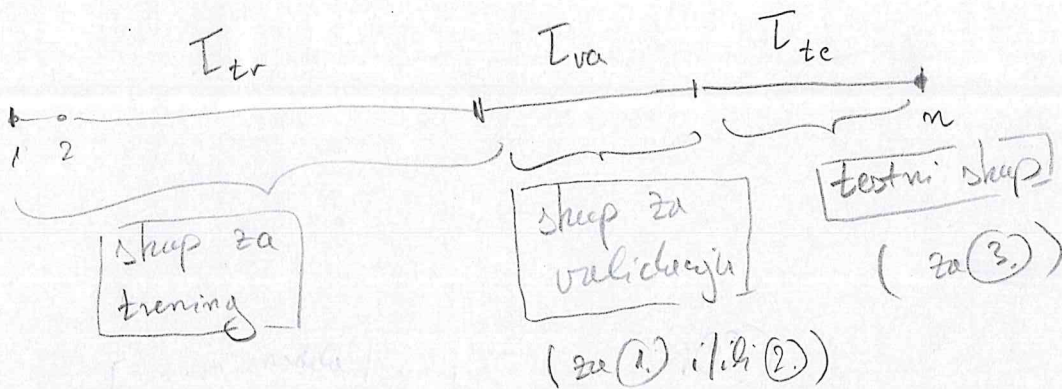
↳ hiperparametar koji kontrolira kompleksnost modela (npr.  $\alpha = \frac{1}{n}$  ili u k-NN metod:).

Ciljevi:

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

Ako je  $n$  velik, najbolji pristup je (na slučajnu način)

podijeliti  $T$ :



1.  $T_{tr}$  → prilagodba raznih modela i/ili istog modela za različite  $\alpha$  [?]
2.  $T_{va}$  → za svaki model ili  $\alpha$  procjenjuje se testna greška (2.1) te npr. odabiremo onaj s najmanjom procijenjenom greškom

3.  $T_{te} \rightarrow$  procenjenjem testnu grešku (2.1) odobrenog modela. 13

Dop. 2.11 [Zašto trebamo  $T_{te}$ ?] (!)

npr. ako je  $T_{ra}$  realizacija od njeke varijabla

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

za  $(X, Y)$ , za nekakvu funkciju  $f$ -ju  $f: \mathbb{R}^p \rightarrow \mathbb{R}$ ,

$$L_{T_{ra}}(f) = \frac{1}{m} \sum_{i=1}^m L(\tilde{Y}_i, f(\tilde{X}^{(i)})) \quad (2.2)$$

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

[Ključna varijabla]  $\forall f!$

ako je  $m$  dovoljno velik.

pa, ako [načinom]  $L_{T_{ra}}(\hat{f}_\pm)$  za sve  $f \in \mathcal{F}$

$$\hat{f} = \underset{f}{\text{argmin}} L_{T_{ra}}(\hat{f}_\pm), \quad (2.3)$$

imamo

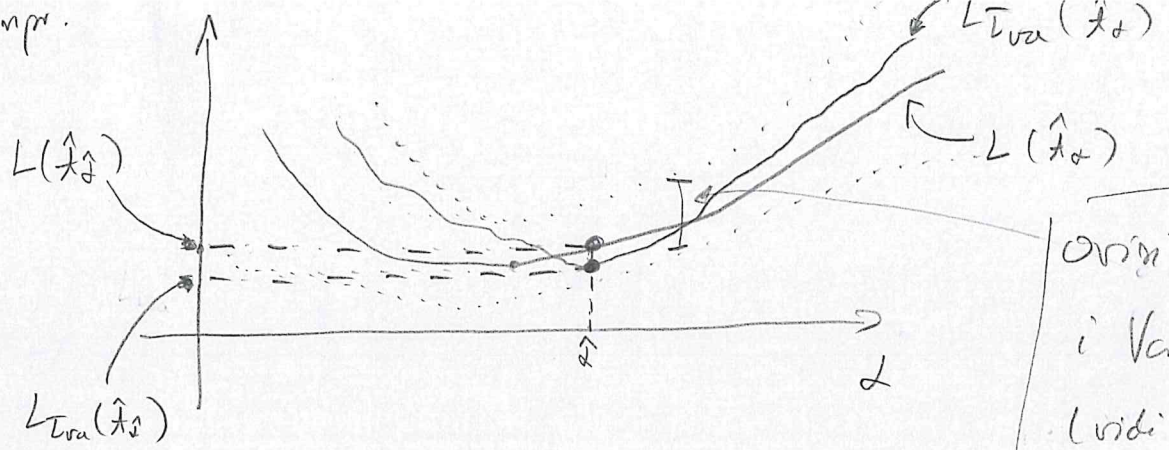
$$E[L_{T_{ra}}(\hat{f}_\pm)] \leq \min_f E[L_{T_{ra}}(\hat{f}_\pm)] = L(\hat{f}_\pm)$$

$\leq L_{T_{ra}}(\hat{f}_\pm), \forall f$  po (2.3)

$\leq L(\hat{f}_\pm)$  [Ključna varijabla]

$$E[L(\hat{f}_\pm) - L_{T_{ra}}(\hat{f}_\pm)] \geq 0! \text{ [probnaost]}$$

npr.



$\text{ovrsti } \underline{m}$   
 $i \text{ Var}(L(Y, \hat{\lambda}_d(x)))$   
 (vidi (2.2))

[veći  $m$ , manja greška]

Šlión, ako imem više velikih modela  $\rightarrow \hat{\lambda}^{(1)}, \dots, \hat{\lambda}^{(k)}$  te

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

imemo

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

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

2.1) Unakrsna validacija (CV)

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

[u tom slučaju] CV [predstavja] najjednostavniju i najpopularniju metodu.

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

15

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

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

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

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

- $\forall j=1, \dots, k$ , računamo [Cross algoritmu]

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

pri čemu je

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

za

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

oni osim  $\mathcal{T}_j$

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

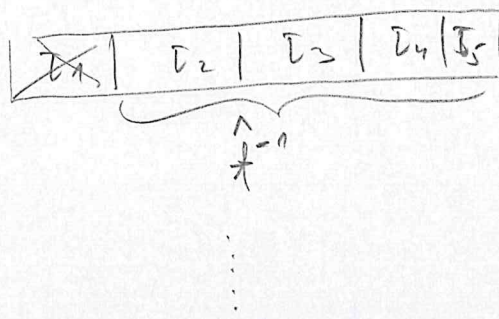
$$CV^{(n)} = L_{cv}^{(n)}(\hat{A}) := \frac{1}{k} \sum_{j=1}^k CV_j^{(n)} \quad (2.5)$$

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

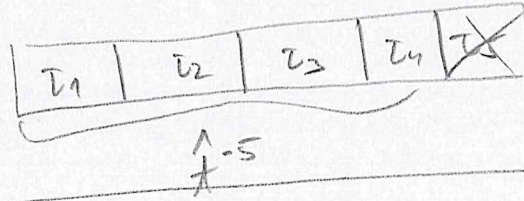
pri čemu je  $j(i) \in \{1, \dots, k\}$  t.d.  $(x^{(i)}, y_i) \in \mathcal{T}_{j(i)}$ .

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

mpm.  $h=5$



$\tau_{va} = \tau_n$   $CV_1^{(5)}$



$\tau_{va} = \tau_5$   $CV_5^{(5)}$

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

↳ u  $j$ -tom koraku  $\tau_j$  igra ulogu skupa za validaciju (tj. merionog testnog skupa).

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

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

Ipak, tipično  $L_{CV}^{(h)}$  koristimo za odabir modela / hiperparametara.

[tipično izabiremo model / hiperparameter koji ima manju testnu grešku  $L(\hat{\tau})$ . [CV je još uvijek težišni nedovoljno istražen!]]

Kako odabrati  $h$ ?

[odlično "skupa za trening"]

$|\tau^{-s}| = m - r_n = \frac{h-1}{h} \cdot m < m$

$\Rightarrow$  tipično imamo  $E_T[L_{CV}^{(h)}(\hat{\tau})] > E_T[L(\hat{\tau})]$

za male  $h$  [odnos velike  $r_n$ ]

pristrenost!

• za velike  $k$ ,  $T^{-1/2}$ ,  $j=1, \dots, k$  su velo bliži, te je

tipično  $\text{Var}_T(L_{CV}^{(h)}(\hat{t}))$  veća  $\rightarrow$  opet BVT!

Nap:!  $\text{Var}_T(L_{CV}^{(h)}(\hat{t}))$  je komplicirana kombinacija

• varijance prosjeka od  $CV_1^{(h)}, \dots, CV_k^{(h)}$

• kovarijance među prognozama  $CV_1^{(h)}, \dots, CV_k^{(h)}$

• u praksi najčešće

(a)  $k=5$  ili  $k=10 \rightarrow T_{tr}^{CV}$  je 80% ili 90% od  $T_{tr}$ .

(b)  $k=n$  za "mali"  $n$  ili ako efikasnije mjerimo izračunati  $L_{CV}^{(h)}(\hat{t})$

$\rightarrow$  tzv. "leave-one-out CV" (LOOCV)  
unijesto LOOCV, bolje ipak  $k < n$  (ali  $k$  velik)

Nap. 2.3 Unijesto  $\hat{t} := t_{min} := \underset{t}{\text{argmin}} L_{CV}^{(h)}(\hat{t})$ , mjerimo uistinu

$\hat{t}$  koji odgovara najjednostavnijem modelu t.d.

$$L_{CV}^{(h)}(\hat{t}_{\hat{t}}) \approx L_{CV}^{(h)}(\hat{t}_{t_{min}}) + \widehat{SE}(t_{min})$$

gdje je

$$\widehat{SE}(t) := \sqrt{\frac{\widehat{\text{Var}}(CV_1^{(h)}, \dots, CV_k^{(h)})}{k}} \quad (2.7)$$

(najini!) prognetelj za stand. grešku od  $L_{CV}^{(h)}(\hat{t}_{\hat{t}})$ .

• "one SE rule"

• problem  $\rightarrow$  (2.7) je što su  $CV_1^{(h)}, \dots, CV_k^{(h)}$  korelirane

- analogno pravilo možemo primijeniti kod odabira modela
- Ozcarov princip: između jednako preciznih modela

13

(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  $\forall \epsilon \in (0, 1/3)$ . Nis algoritom je

1. Na temelju  $T$  izračunaj varijance 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 izlučiti CV algoritmom sa procjenom cjelosti?

→  $i$  (1.) mora biti uključen u sržni korak CV algoritma!

(vidi ESL, poglavlje 7.10.2)

Npr., neka je

•  $X_1, \dots, X_p$  njeđ  $\sim N(0, 1)$  nezavisnom ( $p = 5000$ )

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

⇒  $\forall \hat{\lambda}$ :  $\mathbb{R}^p \rightarrow \text{dom} \hat{\lambda}$  vrijedi

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

10

prop. 0-1 gubitak

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

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

mez. od  $Y$  i  $X$

[maksimum!]

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

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

$\rightarrow$  CV-primer. R

Npr. Ukoliko rad kodiranjem medimor neke metode neodređenog učenja (npr. analiza glavnih komponenti), tj. ne konstantno  $y_i$ -eve, taj koncept nije potrebno uvrstiti u CV algoritam!

## 2.2 Stepnjevi slobode (degrees of freedom) $\rightarrow$ (dt)

$\hookrightarrow$  želim upoređivati kompleksnost različitih metoda

Prop. 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.  $\mathbb{E}[\epsilon_i] = 0$ ,  $\text{Cor}(\epsilon_i, \epsilon_j) = 0$ ,  $\forall i \neq j$

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



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

~~fixed-design nije ovo~~  
 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  
 $E = (\epsilon_1, \dots, \epsilon_n)^T \in \mathbb{R}^{n \times 1}$ , (2.8) ~~postoje~~  
 $Y = \mu + E$ ,  $E[E] = 0$ ,  $Cov(E) = \sigma^2 I_n$   
 korijepijka matrica

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

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

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

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

[dakle, ako je "više" prilagodstvenih podacima, Cov pa samim time i d\hat{A} raste!]

Motivacije za (2.9): usporedit ćemo

$$E \left[ \underbrace{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}_{L_T(\hat{A})} \right]$$

očekivanu grešku na skupu za trening

$$E \left[ E \left[ \underbrace{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}_{L_T(\hat{A})} \right] \right]$$

tzv. očekivana "in-sample" testna greška

gdje je  $\tilde{Y}_i = \{ (x^{(i)}, \tilde{Y}_i) : i=1, \dots, n \}$  testni skup nezavisan 21  
 od  $T(\tilde{Y}_i = \mu_i + \tilde{\epsilon}_i, i=1, \dots, n)$ ,  $(\tilde{\epsilon}_i)_{i=1, \dots, n}$  nez. od  $(\epsilon_i)_{i=1, \dots, n}$ .

Pr. 2.5 | Vredj:

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

tzv. očekivani "optimum"  
 od  $L_T(\hat{\lambda})$

Dokaz |

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

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

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

Pr. 2.6 |

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

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

(iii) Kod lin. regresije metodom najmanjih kvadrata,  $d\lambda(A) = \boxed{p}$  (kasnije)

(iv) Kod kNN regresije,  $d\lambda(A) = \frac{n}{k}$  (DZ)

Svi primjeri iz Pr. 2.6 su oblika

22

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

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. 2.7 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. 2.8 Iz (2.10) slijedi da je statistika

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

nepristran procjenitelj za  $\mathbb{E}[L_T^2(\hat{A})]$

$\rightarrow T$  možemo koristiti za odabir modela.

Pr. 2.2  $\rightarrow$  Kod lin. regresije metodom najmanjih kvadrata,

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

(koristi se za odabir varijabli)