

# 1) Uvod

## 1.1 Problem (mašinski) učenja

Pretp. da imamo podatke

$$T = \{ (x^{(i)}, y_i) : i=1, \dots, m \}$$

uzorak ili "skup za učenje"  
[training set]

gdje je mpr.

- $y_i$  = cijena  $i$ -te nekretnine

odgovor ili zavisna varijabla  
[response]

- $x^{(i)} = (x_{i1}, x_{i2}, \dots, x_{ip})$

↑  
veličina  
(u m<sup>2</sup>)

↑  
pozled  
na  
njeju (DA/NE)

↑  
karakteristike ili  
nezavisne varijable  
(p ∈ N je broj karakteristika)

Najčešće, cilj učenja je predikcija:

$$T \longrightarrow \hat{f} = \hat{f}(T),$$

t.d. je  $\hat{f}(\tilde{x})$  "dobna" procjena za  $\tilde{y}$ , pri čemu je  $(\tilde{x}, \tilde{y})$  "novi" par.

Ako je

- $y$  kvantitativna (neprekidna) varijabla

→ problem regresije (odnaka  $y \in \mathbb{R}$ )

•  $y$  kategorijalna varijabla

→ problem klasifikacije

(BSPROF  $y \in S = \{0, 1, \dots, K-1\} \subseteq \mathbb{R}$ )

1.2 Teorija statističkog učenja

[ statistički pristup problemu učenja  
→ TSU (srednjo) ]

Predpostavke

- odčit  $Y$  je slučajna varijabla

- kovarijate  $X = (X_1, \dots, X_p)$  su slučajni vektor u  $\mathbb{R}^p$

(  $X$  i  $Y$  su zavisni ! npr.  $Y = f(X) + \varepsilon$  )

-  $T$  je jedna realizacija od

$$T = \{ (X^{(i)}, Y_i) : i = 1, 2, \dots, m \},$$

gdje su, tipično,  $(X^{(i)}, Y_i)$ ,  $i = 1, \dots, m$ , mjdl

te  $(X^{(i)}, Y_i) \sim (X, Y)$

Cilj: Za denu  $f$  gubitka  $L: \mathbb{R}^2 \rightarrow \mathbb{R}$  [ loss function ]  
tražimo  $\hat{f} = \hat{f}(T): \mathbb{R}^p \rightarrow \mathbb{R}$  sa što manjom bestnom

gneskom

$$L(\hat{f}) := \mathbb{E}[L(Y, \hat{f}(X))] \quad (1.1)$$

[  $L(y, \hat{f}(x))$  je "trošček" predikcije  $\hat{f}(x)$  ako je strama ujednost  $y$  ] [ 3

Nap. 1 Ako imamo realizaciju

$$\tilde{\mathcal{D}} = \{ (\tilde{x}^{(i)}, \tilde{y}_i) : i=1, \dots, m \}$$

testni skup

jes jechovy mnojot wawoku iz  $(x, y)$ , za veliki  $m$ ,

$$L_{\tilde{\mathcal{D}}}(\hat{f}) := \frac{1}{m} \sum_{i=1}^m L(\tilde{y}_i, \hat{f}(\tilde{x}^{(i)})) \approx L(\hat{f}). \quad (1.2)$$

greška na  $\tilde{\mathcal{D}}$

Primeri izbori za  $L$  su npr.

•  $\forall y \in \mathbb{R} \rightarrow L(y, \hat{f}(x)) = (y - \hat{f}(x))^2 \quad (1.3)$

•  $\forall y \in \{0, 1\} \rightarrow L(y, \hat{f}(x)) = \begin{cases} 1, & y \neq \hat{f}(x) \\ 0, & y = \hat{f}(x) \end{cases} \quad (1.4)$

Uz oimw,  $\forall \hat{f}: \mathbb{R}^P \rightarrow \mathbb{R}$ , ako je za  $x \in \mathbb{R}^P$  defin.

$$L_x(\hat{f}) := E[L(Y, \hat{f}(x)) | X=x], \quad (1.5)$$

ujedn.

testna greška od  $\hat{f}$  u  $x$

$$L(\hat{f}) = \int_{\mathbb{R}^P} L_x(\hat{f}) \mathbb{P}_X(dx). \quad (1.6)$$

Nap. 1

(i)  $\mathbb{P}_X(B) := \mathbb{P}(X \in B), \forall B \in \mathbb{R}^P$  Borelov  
[ distribucija od  $X$  ]

100)

$$\mathbb{E}[L(Y, t(x)) | X=x] := h(x),$$

4

pri čemu je  $h: \mathbb{R}^p \rightarrow \mathbb{R}$  t.o.

$$h(X) = \mathbb{E}[L(Y, t(X)) | X] \quad \text{g.o.}$$

memo

$$\mathbb{E}[L(Y, t(x))] = \mathbb{E}[\mathbb{E}[L(Y, t(x)) | X]]$$

$$= \mathbb{E}[h(X)] = \int_{\mathbb{R}^p} h(x) P_X(dx)$$

Iz (1.6) sledi da  $t^*: \mathbb{R}^p \rightarrow \mathbb{R}$  defin. sa

$$t^*(x) := \underset{t \in \mathbb{R}}{\text{argmin}} \mathbb{E}[L(Y, t) | X=x], \quad x \in \mathbb{R}^p \quad (1.7)$$

minimizira (1.1), tj:  $L_x(t^*) \leq L_x(t)$  i  $L(t^*) \leq L(t)$ ,

za sve  $t$  i sve  $x \in \mathbb{R}^p$ .

optimalna ("reducirana") greska

Pr. 1

(i) Ako je  $L$  iz (1.3),

$$t^*(x) = \mathbb{E}[Y | X=x], \quad (1.8)$$

(ii) Ako je  $L$  iz (1.4),

$$t^*(x) = \begin{cases} 1, & \text{ako } P(Y=1 | X=x) > \frac{1}{2}, \\ 0, & \text{inače.} \end{cases} \quad (1.9)$$

Bayesov klasifikator

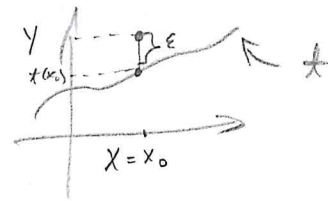
(DZ - tretinajite  $P(\cdot | X=x)$ ;  $E[\cdot | X=x]$  kao standardne  $P(\cdot)$ ;  $E[\cdot]$ ) [5]

[  $f^*$  merimo ne znamo jer ne znamo distribuciju od  $(X, Y)$  ]

[  $\hat{f}$  binarno iz neke familije  $f$ -ja  $\mathcal{F}$  koristeći neki algoritam  $\rightarrow$  PAC učenje, VC dimenzija... ] (TSU)

1.3 Primer

Pretpostavimo model



$$Y = f(x) + \epsilon, \quad (1.10)$$

gdje (i)  $f$  neka  $f$ -ja, (ii)  $\epsilon$  sl. varijabla meraciona od  $X$ , (iii)  $E[\epsilon] = 0$ ,  $\text{Var}(\epsilon) = \sigma^2 < \infty$ .

Koćino, nužno je

$$\begin{aligned} f^*(x) &= E[Y | X=x] = E[f(x) + \epsilon | X=x] \\ &= f(x) + E[\epsilon | X=x] = f(x) + \underbrace{E[\epsilon]}_{=0} = f(x). \end{aligned}$$

Nadalje,  $\forall x \in \mathbb{R}^p$  irreducibilna greška je

$$L_x(x) = E[\underbrace{(Y - f(x))}_{\epsilon^2} | X=x] \stackrel{(1.10)}{=} \sigma^2,$$

te

$$L(f) = E[(Y - f(x))^2] = \sigma^2.$$

na to me moramo utjecati:  $\epsilon$  predstavlja onaj dio info od  $Y$  koji nije sadržan u  $x$ .

1.3.1) Metoda k najbližih susjeda (k-nearest neighbors)  $\rightarrow$  kNN [6]

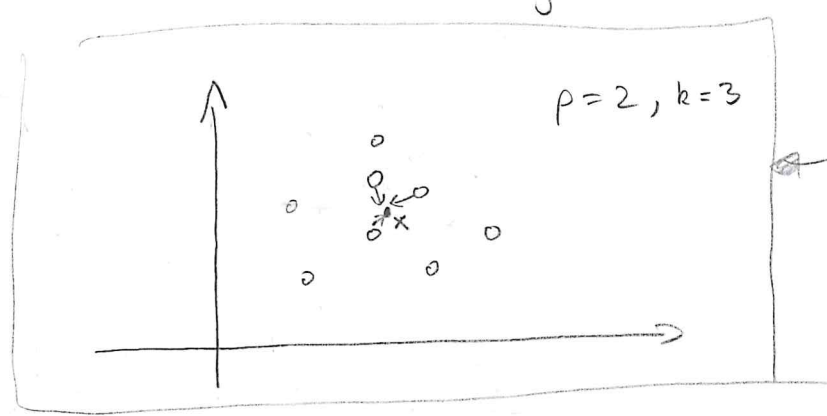
Za  $k \in \{1, \dots, m\}$  fiksiran,

$$\hat{f}_k(x) = \frac{1}{k} \sum_{x^{(i)} \in N_k(x)} y_i, \quad x \in \mathbb{R}^p, \quad (1.11)$$

pri čemu je

$N_k(x)$  = skup od  $k$  točaka iz

$\mathcal{T}_x = \{x^{(1)}, \dots, x^{(m)}\}$  koje su najbliže točki  $x$



$\hat{f}_3(x)$  je prosjek vrijednosti  $y$ -a za 3 najbliža susjeda od  $x$  u  $\mathcal{T}$

Dop.1

(i) neparametarski pristup -  $\hat{f}_k(x)$  direktno procjenjuje

$$f^*(x) = f(x) = E[Y | X=x] \quad \left[ \begin{array}{l} \text{jedini zahtjev je da je} \\ f \text{ dovoljno glatka.} \end{array} \right]$$

(ii) [hiperparametar] parametar  $k$  kontrolira "kompleksnost" metode (tj. "fleksibilnost" od  $\hat{f}_k$ ), npr.

•  $k=1 \Rightarrow \hat{f}_1(x^{(i)}) = y_i, \forall i=1, \dots, n!$   
 ("overfitting" jer  $L_{\mathcal{T}}(\hat{f}_1) = 0$ )

•  $k=m \Rightarrow \hat{f}_m(x) = \frac{1}{n} \sum_{i=1}^m y_i = \bar{y}, \forall x \in \mathbb{R}^p!$   
 $N_m(x) = \mathcal{T}_x$  ("underfitting")

(iii) U praksi,  $k$  odabiremo unakrsnom validacijom.  
(kasnije)

7

### 1.3.2 Linearna regresija

Pretpostavljamo da je zve neki  $\beta = (\beta_0, \dots, \beta_p) \in \mathbb{R}^{p+1}$ ,

$$f(x) = \mathbb{E}[Y | X=x] = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p, \quad (1.12)$$

$=: f_{\beta}(x)$

za sve  $x = (x_1, \dots, x_p) \in \mathbb{R}^p$ , a  $\beta$  procjenjujemo npr. metodom  
najmanjih kvadrata  $\rightarrow \hat{\beta} = (\hat{\beta}_0, \dots, \hat{\beta}_p)$ , te

stavimo

$$\hat{f}(x) = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p, \quad x \in \mathbb{R}^p.$$

Nap. 1

(i) parametarski pristup - pretpostaviti smo da  
je  $f \in \mathcal{F}_{\beta} = \{ f_{\beta} : \beta \in \mathbb{R}^{p+1} \}$  pa čemu  
je broj parametara  $(p+1)$  končan.

(ii) metoda ima smisla: kada (1.12) ni  
približno ne vrijedi  $\rightarrow$  aproksimiramo  $f$   
sa  $f$ -jom iz  $\mathcal{F}_{\beta}$ .

$\rightarrow$  uvodno - predavanje. pdt  
(knn vs. lin. regresija)

1.4) Odnos između pristranosti i varijance (!)<sup>L</sup>  
 (bias - variance tradeoff) → (BVT)

[najvažnija stvar kada govorimo o predikciji!]

[terminologija koju ćemo koristiti stalno]

Pretp. da je  $kov$  u (1.10),

$$Y = f(X) + \epsilon.$$

Ako posmatramo  $\hat{f} = \hat{f}(T)$  gdje je  $T$  slučajni uzorak,  $\hat{f}$  postaje slučajna  $f$ -ja,

a  $L_x(\hat{f})$  i  $L(\hat{f})$  slučajne varijable.

Prop. 1.11 (BVT)

Ako je  $L(y, \hat{y}) = (y - \hat{y})^2$ ,  $\forall x \in \mathbb{R}^d$  ujed.

$$\mathbb{E}_T[L_x(\hat{f})] = \sigma^2 + \text{Var}_T(\hat{f}(x)) + (\mathbb{E}_T[\hat{f}(x)] - f(x))^2 \quad (1.13)$$

↑ ↑ ↑  
 ireducibilna greška + varijanca + pristranost<sup>2</sup> (bias)

"očekivana testna greška u  $x$ "

"reducibilna greška" [odj je to smanjiti]

Nap. 1  
 (i)  $\mathbb{E}_T$  predstavlja očekivanje u kojem slučajnost dolazi od  $T$ .

(ii) (1.13)  $\Rightarrow \mathbb{E}_T[L(\hat{f})] = \int_{\mathbb{R}^d} \mathbb{E}_T[L_x(\hat{f})] P_X(dx)$  (1.6) + FUBINISEU TH.



10)  $\forall g: \mathbb{R}^p \rightarrow \mathbb{R}; \forall x \in \mathbb{R}^p,$

$$\begin{aligned}
 L_x(g) &= \mathbb{E}[(Y - g(X))^2 \mid X=x] = \mathbb{E}[(f(x) - g(x) + \varepsilon)^2 \mid X=x] \\
 &= \mathbb{E}[(f(x) - g(x))^2 + (f(x) - g(x))\varepsilon + \varepsilon^2 \mid X=x] \\
 &= (f(x) - g(x))^2 + (f(x) - g(x)) \underbrace{\mathbb{E}[\varepsilon \mid X=x]}_{= \mathbb{E}[\varepsilon] = 0} + \underbrace{\mathbb{E}[\varepsilon^2 \mid X=x]}_{= \mathbb{E}[\varepsilon^2] = \sigma^2} \\
 &= (f(x) - g(x))^2 + \sigma^2
 \end{aligned}$$

20)  $\forall x \in \mathbb{R}^p,$

$$\mathbb{E}_T[L_x(\hat{f})] \stackrel{10)}{=} \sigma^2 + \mathbb{E}_T[(f(x) - \hat{f}(x))^2]$$

$\downarrow$   
 ozn. d T  
 $\pm \mathbb{E}_T[\hat{f}(x)]$

$$\begin{aligned}
 &= \sigma^2 + \mathbb{E}_T[(f(x) - \mathbb{E}_T[\hat{f}(x)])^2] \\
 &\quad + \underbrace{(f(x) - \mathbb{E}_T[\hat{f}(x)])}_{\text{konstanta}} \underbrace{\mathbb{E}_T[\mathbb{E}_T[\hat{f}(x)] - \hat{f}(x)]}_{= 0} \\
 &\quad + \mathbb{E}_T[(\mathbb{E}_T[\hat{f}(x)] - \hat{f}(x))^2]
 \end{aligned}$$

$$= \sigma^2 + \underbrace{(f(x) - \mathbb{E}_T[\hat{f}(x)])^2}_{\text{pristrenost}^2} + \text{Var}_T(\hat{f}(x))$$

Tipično, kako kompleksnost metode raste, pristranost  $\hat{f}$  se smanjuje, a varijanca povećava.

Pr. 1 u kNN metodi,

(i) ako je  $k=1$ ,  $\forall x \in \mathbb{R}^D$ ,

$$\hat{f}(x) = f(X^{(i)}) + \epsilon \quad (1.14)$$

→ za veliki  $n$  (pod nekim uvjetima),

$$\hat{f}(X^{(i)}) \approx f(x), \text{ g.d.} \quad (1.15)$$

te

$$E_T [\hat{f}(x)] \approx f(x) + E[\epsilon] = f(x)$$

i

$$\text{Var}_T [\hat{f}(x)] \approx 0 + \text{Var}(\epsilon) = \sigma^2$$

⇒ ako je  $n$  velik i  $\sigma$  ne prevelik, imamo

mala pristranost i veliku varijancu

↑ (dodaj: od  $\epsilon$ -a u (1.14))

(ii) ako je  $k=m$ ,  $\forall x \in \mathbb{R}^D$ ,

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n f(X^{(i)}) + \frac{1}{n} \sum_{i=1}^n \epsilon_i,$$

g.d. za veliki  $n$ ,

→ za veliki  $n$ ,

$$\hat{f}(x) \approx E[f(X)] + E[\epsilon]$$

$$= E[f(X)] \quad \text{g.d.}$$

ista konstanta  $\forall x \in \mathbb{R}^D$ !

te

$$E_T [\hat{f}(x)] \approx E[f(X)],$$

i

$$\text{Var}_T [\hat{f}(x)] \approx 0 \quad [\text{ne heuristički!}]$$

$\Rightarrow$  Za veliki  $n$ , ako  $f$  nije konstantna  $f$ -je, imamo (1)

veliku pristranost i malu varijancu

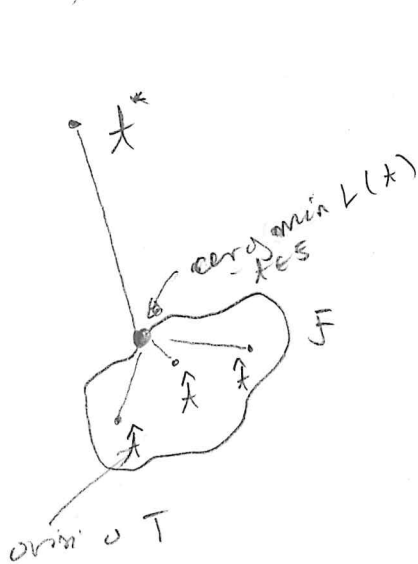
jer  $f(x)$  procjenjujemo  
konstanti: sk  $f(x^{(i)})$

Nop-1 BVT je fundamentalna stvar u stat. učenju  
(kada govorimo o predikciji), i jačlja se  
bez obzira na oblik  $f$ -je  $L$ .

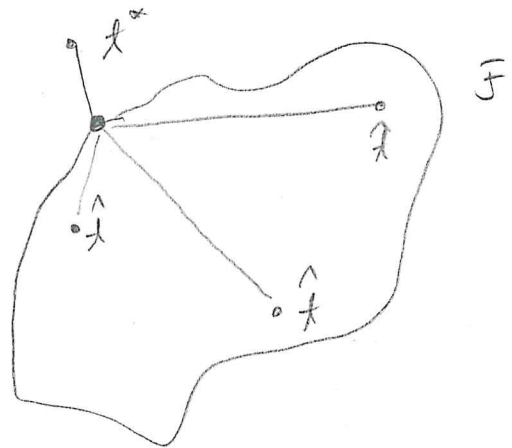
$$L(\hat{f}) = L(f^*) + (L(\hat{f}) - L(f^*))$$

$$= \underbrace{L(f^*)}_{\text{ireducibilna greška}} + \underbrace{(L(\hat{f}) - \inf_{f \in \mathcal{F}} L(f))}_{\text{greška procjene ("varijanca")}} + \underbrace{(\inf_{f \in \mathcal{F}} L(f) - L(f^*))}_{\text{greška aproksimacije ("pristranost")}}$$

Intuitivno,



velika pristranost  
mala varijanca



mala pristranost  
velika varijanca

$\rightarrow$  uvodno predavanje pet (Odnos između pristranosti i varijance.)