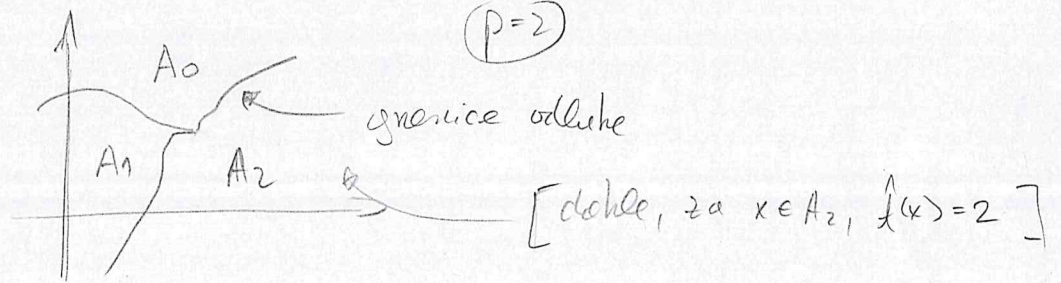


(7) Osnovne metode za klasifikaciju → (7.1) Bayesov klasifikator

↳ $Y \in S = \{0, 1, \dots, K-1\}$ za $K \geq 2$, $X \in \mathbb{R}^P$ [70]

⇒ tražimo $\hat{f} = \hat{f}(x): \mathbb{R}^P \rightarrow S$

Svrha \hat{f} dijeli (\mathbb{R}^P) na $A_k := \hat{f}^{-1}(\{k\})$, $k \in S$



Svrha \hat{f} -ja gubitka $L: S \times S \rightarrow \mathbb{R}$ je zapravo $K \times K$ matrica

[$L(a,b)$ = kazna ako stvaramo klasu a predviđamo b]

↳ $L(\hat{f}) := \mathbb{E}[L(Y, \hat{f}(X))]$, te $\forall x \in \mathbb{R}^P$

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

$$= \sum_{k \in S} L(k, \hat{f}(x)) \underbrace{\mathbb{P}(Y=k | X=x)}_{=: p_k(x)} \quad (7.1)$$

↳ Bayesov klasifikator je $\hat{f}^*: \mathbb{R}^P \rightarrow S$ dan s

$$\hat{f}^*(x) = \underset{c \in S}{\operatorname{argmin}} \sum_{k \in S} L(k, c) p_k(x), \quad (7.2)$$

te očito vrijedi $\forall \hat{f}: \mathbb{R}^P \rightarrow S$, $\forall x \in \mathbb{R}^P$,

$$L_x(\hat{f}^*) \leq L_x(\hat{f}), \quad \boxed{L(\hat{f}^*) \leq L(\hat{f})}$$

[Naravno, \hat{f}^* ne znamo
jer ne znamo
 $p_k(x)$]

"Bayesov način"

7.2 Logistička regresija

72

za $(K=2)$, tj. $S = \{0, 1\}$, to je GLM za $\{B(1, p)\}$ familiju uz kanonsku t-ju veze g (tj. logit t-ja): za $\beta \in \mathbb{R}^p$ pretpost.

$$\log \left(\frac{p_1(x)}{1 - p_1(x)} \right) = g(p_1(x)) \equiv x^T \beta, \quad x \in \mathbb{R}^p \quad (7.6)$$

tj.

$$p_1(x) = g^{-1}(x^T \beta) \equiv \frac{e^{x^T \beta}}{1 + e^{x^T \beta}}, \quad x \in \mathbb{R}^p \quad (7.7)$$

$\frac{p_1(x)}{p_0(x)}$ su "odds" (engl. odds), npr.

• $p_1 = \frac{2}{3} \Rightarrow p_0 = \frac{1}{3}$; odds = $\frac{2}{1}$

• $p_1 = \frac{2}{5} \Rightarrow p_0 = \frac{3}{5}$; odds = $\frac{2}{3}$

Interpretacija parametara β ?

uz (7.6) (tj. (7.7)), za $x = (x_1, \dots, x_p) \in \mathbb{R}^p$ te

$$x' = (x_1, \dots, x_{j-1}, \overbrace{x_j + 1}, x_{j+1}, \dots, x_p) \in \mathbb{R}^p, \text{ imamo}$$

$$\log \left(\frac{p_1(x')}{p_0(x')} \right) = \log \left(\frac{p_1(x)}{p_0(x)} \right) + \beta_j, \text{ tj.}$$

$$\frac{p_1(x')}{p_0(x')} = \frac{p_1(x)}{p_0(x)} \cdot \underbrace{e^{\beta_j}}_{\begin{matrix} > 1, \text{ ako } \beta_j > 0 \\ < 1, \text{ ako } \beta_j < 0 \end{matrix}}$$

~~powerje, R~~

7.2.1 Multinomijalna logistička regresija ($K \in \mathbb{N}$)

pretp. da za neke $\beta^{(1)}, \dots, \beta^{(K-1)} \in \mathbb{R}^p$,

$$\log \left(\frac{p_k(x)}{p_0(x)} \right) = x^T \beta^{(k)}, \quad x \in \mathbb{R}^p, \quad k=1, \dots, K-1 \quad (7.8)$$

[$K-1$ nezavisnih "log. regr." za (K) u odnosu na (0) .]

Uopzi kategorija \mathcal{O} je tzv. bazno katezija (engl. baseline).
 prirodno!

(7.8) \Rightarrow za x, x' kao gore ($x' = x + e_j$),
 $\forall j = 1, \dots, P$

$$\frac{p_k(x')}{p_0(x')} = \frac{p_k(x)}{p_0(x)} \cdot e^{\beta_j^{(k)}}$$

73

(DZ) Pokažite da je (7.8) ekvivalentno s

$$p_k(x) = \frac{e^{x^T \beta^{(k)}}}{1 + \sum_{i=1}^{K-1} e^{x^T \beta^{(i)}}}, \quad k=1, \dots, K-1, \quad (7.9)$$

$$p_0(x) = \frac{1}{1 + \sum_{i=1}^{K-1} e^{x^T \beta^{(i)}}}$$

ako su procjene $\hat{\beta}^{(1)}, \dots, \hat{\beta}^{(K-1)}$ te $\hat{\beta}^{(0)} := (0, \dots, 0) \in \mathbb{R}^P$,

$$\hat{p}_k(x) := \frac{e^{x^T \hat{\beta}^{(k)}}}{1 + \sum_{i=0}^{K-1} e^{x^T \hat{\beta}^{(i)}}}, \quad k=0, 1, \dots, K-1. \quad (= \text{softmax}_k(x^T \hat{\beta}^{(0)}, \dots, x^T \hat{\beta}^{(K-1)}))$$

[β_j : 0-1 gukitech]

ako je $\hat{f}(x) = \arg \max_{k \in S} \hat{p}_k(x)$, imamo

$$\hat{f}(x) = \arg \max_{k \in S} \underbrace{x^T \hat{\beta}^{(k)}}_{=: \hat{\delta}_k(x)}, \quad x \in \mathbb{R}^P$$

"diskriminacijske f -je"

Granice odluke između klasa $k, l \in S$ je

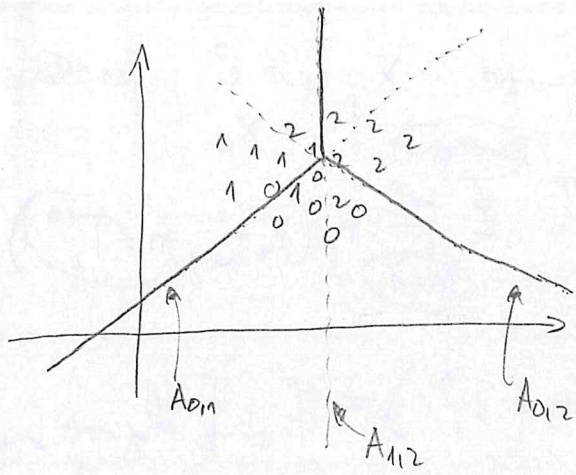
$$A_{k,l} = \{x \in \mathbb{R}^P : \hat{\delta}_k(x) = \hat{\delta}_l(x)\} = \{x : x^T (\hat{\beta}^{(k)} - \hat{\beta}^{(l)}) = 0\}$$

= hiperravna u \mathbb{R}^P , $\forall k, l$

linearna granica odluke!

[tipično imamo i $x_0 = 1$ pa $A_{k,l}$ ne prolazi nužno kroz 0 .]

npv.1



$K=3, P=2$

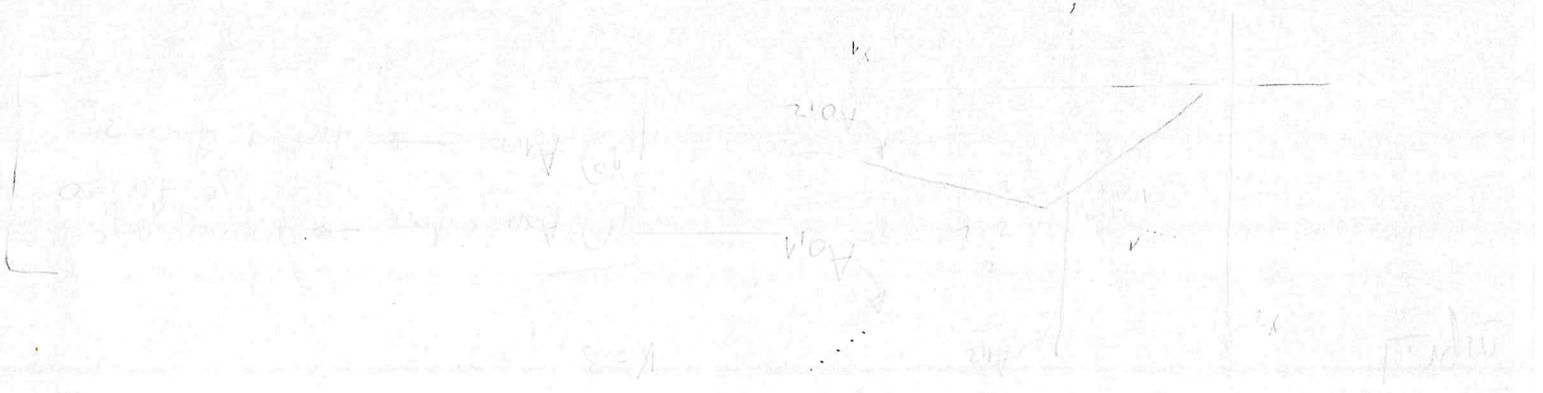
[74]

(1^o) A_{011} i A_{012}
 → merketaj $\lambda x: f(x)=03$

(2^o) $A_{112} \rightarrow \lambda f(x)=13$ i $f(x)=23$

[Koliko li se A_{112} spoci u jednoj tocki? PA]

Uop.! Multivarijantna log. regresija je zapravo specijalan slučaj tzv. "feedforward" neuralne mreže koja nema "skrivenih" slojeva. ■



7.3) Generativni modeli

• u log. regrezi direktno modeliramo $p_k(x)$, $k \in S$
→ [primer] "diskriminacijska" metoda

• alternativno, ako modeliramo $f_k(x) := P(X=x | Y=k)$,
 $x \in \mathbb{R}^p$, $k \in S$, te $\pi_k := P(Y=k)$,

$$p_k(x) = \frac{f_k(x) \pi_k}{\sum_{i \in S} f_i(x) \pi_i} \quad (7.10)$$

Bayes

• Zapravo modeliramo $P(X=x, Y=k) = f_k(x) \pi_k$, tj.
 $P((X, Y) \in \cdot) \rightarrow$ "generativne" metode

- 1) Linearna diskriminacijska analiza (LDA)
- 2) Koadnutna " " " " (QDA)
- 3) Naivni Bayes

Uočimo, (7.10) povlači (uz 0-1 gubitak)

$$f^*(x) = \operatorname{argmax}_{k \in S} f_k(x) \pi_k \quad (7.11)$$

7.3.1 LDA

• Pretpostavka je da $X | Y=k \sim N(\mu_k, \Sigma_k)$, tj.

$$f_k(x) = \frac{1}{(2\pi)^{\frac{p}{2}} \sqrt{\det(\Sigma_k)}} \cdot \exp\left\{-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1}(x-\mu_k)\right\}, x \in \mathbb{R}^p, \quad (7.12)$$

te dodatno

$$\Sigma_k = \Sigma, \forall k \in S! \quad (7.13)$$

$f^*(x) = \underset{k \in S}{\operatorname{argmax}} \log(\pi_k(x) \cdot \pi_k)$
 $= \log(\pi_k(x)) + \log(\pi_k)$

(Mahalevskijeva udaljenost od x i μ_k)

$\stackrel{(7.12)}{=} \underset{k \in S}{\operatorname{argmin}} \left\{ \frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) - \log(\pi_k) \right\}, x \in \mathbb{R}^p$ (7.14)
 $=: d_k(x)$

Ako je $\Sigma = I_p$ te $\pi_0 = \dots = \pi_{k-1}$,

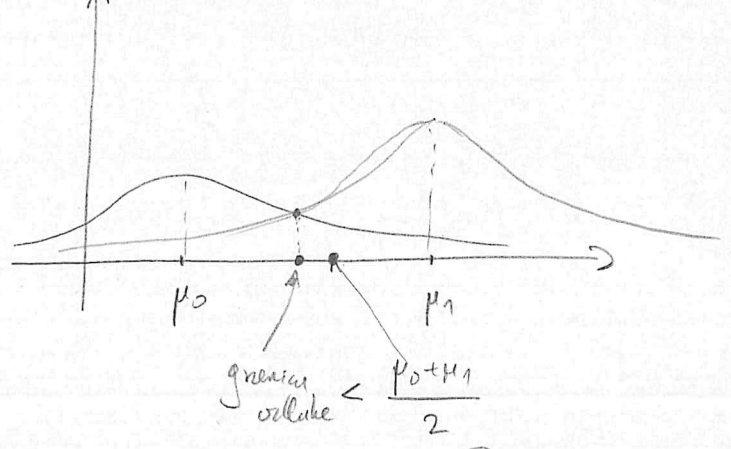
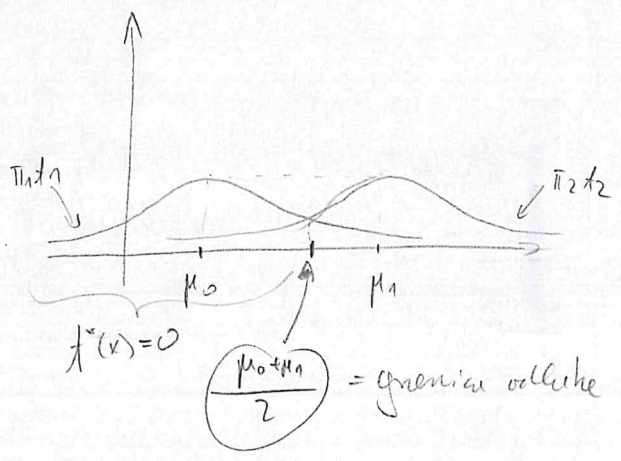
$f^*(x) = \underset{k \in S}{\operatorname{argmin}} \|x - \mu_k\|_2^2$

$f_k = f_k(\mu_k, \sigma^2)$

Pr. | $p=1, k=2$ ($\Sigma = \sigma^2 > 0$)

1° $\pi_0 = \pi_1$

2° $\pi_1 > \pi_2$



[uzimamo u obzir da je vise $(x_i, y_i) \ni \{y_i = 1\}$]

$d_k(x) = \frac{1}{2} x^T \Sigma^{-1} x - x^T \underbrace{\Sigma^{-1} \mu_k}_{=: w_k} + \underbrace{\mu_k^T \Sigma^{-1} \mu_k - \log(\pi_k)}_{=: -b_k}$

\leftarrow ne ovini $\sigma^2(h)$

$f^*(x) = \underset{k \in S}{\operatorname{argmax}} x^T w_k + b_k$

\rightarrow linearna diskriminacijska \rightarrow "LDA"
 f je

Specijalno, LDA \rightarrow linearna granica odluke.

Parametre projekcije su:

$$\hat{\pi}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{y_i = k\}},$$

$$\hat{\mu}_n = \frac{\sum_{i=1}^n x^{(i)} \mathbb{1}_{\{y_i = k\}}}{\sum_{i=1}^n \mathbb{1}_{\{y_i = k\}}},$$

"pooled" varianci kov. matrica

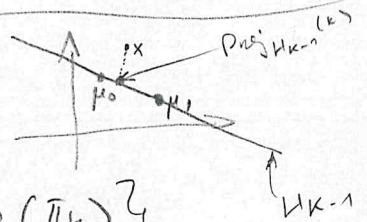
$$\hat{\Sigma} = \frac{1}{n - (k)} \sum_{i=1}^n (x^{(i)} - \hat{\mu}_{y_i})(x^{(i)} - \hat{\mu}_{y_i})^T.$$

dim. $k-1!$

$$\hat{f}_k(x) = \operatorname{argmin}_{h \in S} \hat{d}_k(x).$$

Nap. 1 (i) (ii) Pretp. da je $\hat{\Sigma} = \bar{\Gamma}_p$, te $H_{k-1} := \operatorname{aff}(\mu_0, \dots, \mu_{k-1})$

$$\hat{f}_k(x) = \operatorname{argmin}_{h \in S} \left\{ \frac{1}{2} \|x - \mu_h\|_2^2 - \log(\pi_h) \right\}$$



$$= \operatorname{argmin}_{h \in S} \left\{ \frac{1}{2} \|\operatorname{Proj}_{H_{k-1}}(x) - \mu_h\|_2^2 - \log(\pi_h) \right\}$$

\Rightarrow dim je $k-1 < p$, imamo smaljuje dimenzije

(ii) općenito, prv napomenimo transformaciju

$$x \mapsto x^* := \hat{\Sigma}^{-\frac{1}{2}} x \quad (\Rightarrow (x - \hat{\mu}_h)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_h) = \|x^* - \hat{\mu}_h^*\|_2^2, \hat{\Sigma}^* = \bar{\Gamma}_p)$$

te gledamo projekcije na $H_{k-1}^* = \operatorname{aff}(\mu_0^*, \dots, \mu_{k-1}^*)$

Može se pokušati da postoje vektori $[v_1, \dots, v_{k-1}] \in \mathbb{R}^p$

t.d.

$$\hat{f}_k(x) = \operatorname{argmin}_{h \in S} \left\{ \frac{1}{2} \|V \cdot x - V \cdot \hat{\mu}_h\|_2^2 - \log(\pi_h) \right\}$$

$$\text{za } V = \begin{bmatrix} -v_1^T \\ \vdots \\ -v_{k-1}^T \end{bmatrix} \in \mathbb{R}^{(k-1) \times p}$$

$(v_i^T \cdot x)$ je tzv. i -ta diskriminacijska koordinate

Primjenjena statistika i ESL 4.3.3

\rightarrow Wine-LDA.R

→ iste kao kod LDA, ali bez pretp. $\Sigma_k = \Sigma, \forall k \in S$.

$$\Rightarrow f^*(x) = \underset{k \in S}{\operatorname{argmax}} \left\{ -\frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) - \frac{1}{2} \log(\det(\Sigma_k)) + \log \pi_k \right\}$$

$$= \underset{k \in S}{\operatorname{argmax}} \left\{ -\frac{1}{2} x^T \Sigma_k^{-1} x + \dots \right\}$$

kvadratne granice odluke.

Ipak, u odnosu na LDA, QDA progleduje $\hat{\Sigma}_k, k \in S$

→ dodatnih $\frac{p(p-1)}{2} \cdot (K-1)$ parametara!

Kada je p velik, čak i LDA progleduje previše parametara

→ Σ, μ_k, π_k
 $\frac{p(p-1)}{2} \cdot K \cdot p \quad K-1$

7.3.3 Naivni Bayes (NB)

→ pretp. da su $\forall k \in S, X_1, \dots, X_p$ nezavisne uz da su $Y=k, t.j.$

$$f_k(x) = \prod_{j=1}^p f_{k,j}(x_j), \quad \forall x = (x_1, \dots, x_p) \in \mathbb{R}^p \quad (7.15)$$

[nezavisnost je oblik regularizacije]

Marginalne gustoće mišerni progledati npr.

• ako $X_j \in \mathbb{R} \Rightarrow$ pretpostavimo

$$X_j | Y=k \sim N(\mu_{k,j}, \sigma_{k,j}^2)$$

(kao QDA uz $\Sigma_k =$ dijagonalna matrica)

• ako $x_j \in S_j = \{0, \dots, k_j - 1\} \Rightarrow$ mpr, neprenosivost

$$\hat{\pi}_{k,j}(x) = \frac{\sum_{i=1}^n \mathbb{1}_{\{x_i = k, x_{i,j} = x\}}}{\sum_{i=1}^n \mathbb{1}_{\{x_i = k\}}}$$

7.3.4 Usporedna metoda [dani]

• log. razmjena — $\log\left(\frac{p_k(x)}{p_0(x)}\right) = \boxed{X^T \beta^{(k)}}$ (tipično $x_0 = 1$)

• LDA — $\log\left(\frac{p_k(x)}{p_0(x)}\right) \stackrel{(2)}{=} \log \frac{\pi_k}{\pi_0} = \frac{1}{2} (\mu_k + \mu_0)^T \Sigma^{-1} (\mu_k - \mu_0) + \boxed{X^T} \Sigma^{-1} (\mu_k - \mu_0)$

$=: \alpha_{k,0} + \sum_{i=1}^{k-1} \alpha_{k,i} \cdot x_i$

$=: \boxed{X^T \alpha^{(k)}}$
 \uparrow
 $x_0 = 1$

(i) ipak, $\hat{\alpha}^{(m)} \neq \hat{\beta}^{(m)}$ [te je log. razmjena općenitija jer ne pretpostavlja ništa o razd. $P(X \in \cdot)$].

• NB — $\log\left(\frac{p_k(x)}{p_0(x)}\right) \stackrel{(7.15)}{=} \log\left(\frac{\pi_k}{\pi_0}\right) + \sum_{j=1}^p \log\left(\frac{f_{k,j}(x_j)}{f_{0,j}(x_j)}\right)$

$=: \alpha_{k,0} + \sum_{j=1}^p \alpha_{k,j}(x_j)$

↳ fleksibilne granice odlike

↳ ipak, veća interakcija

• QDA — $\log\left(\frac{p_k(x)}{p_0(x)}\right) = \alpha_k^1 + \sum_{j=1}^p \alpha_{k,j} x_j + \boxed{\sum_{j_1, j_2=1}^p \alpha_{k,j_1 j_2} x_{j_1} x_{j_2}}$

↑
interakcije

7.4 kNN metoda

↳ za $h_0 \in \{1, \dots, n\}$,

$$\hat{p}_h(x) = \frac{1}{n_0} \sum_{x^{(i)} \in M_{h_0}(x)} \mathbb{1}_{\{y_i = h\}}, \quad h \in S, x \in \mathbb{R}^p$$

$$f(x) = \underset{h \in S}{\operatorname{argmax}} \hat{p}_h(x) = \underset{h \in S}{\operatorname{argmax}} \sum_{x^{(i)} \in M_{h_0}(x)} \mathbb{1}_{\{y_i = h\}}$$

Ucap. | za velike p

→ problem dimenzionalnosti!

↑ "majority vote"
[uz 0-1 gubitak]

→ S-and-P.R