

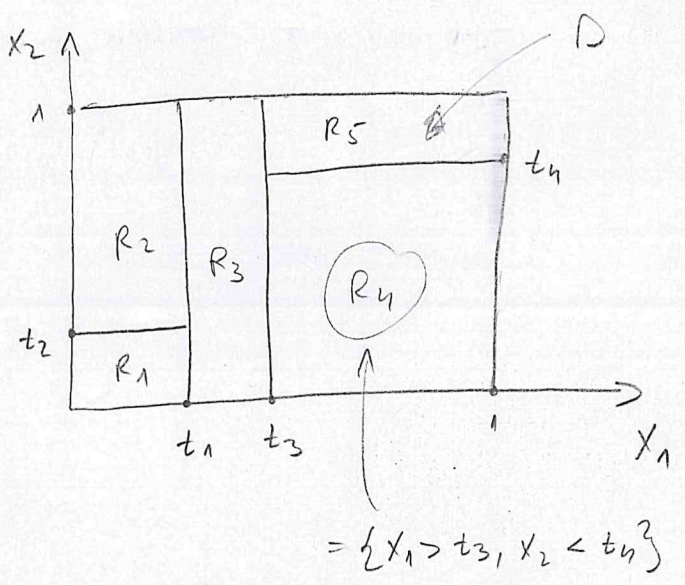
8. Metode bazirane na stablima

8.1 CART (Classification and Regression Trees)

• dijelimo prostor prediktora $D \subseteq \mathbb{R}^p$ na p -dimenz. pravokutnike te prilagodavamo jednostavni model na svakom pravokutniku

• koristimo rekurzivne binarne particije

Pr. 8.1) $(y \in \mathbb{R})$, $X = (x_1, x_2) \in [0,1]^2 =: D$, te x_1 i x_2 kvantitativne.



1. D dijelimo u $x_1 = t_1$
2. $\{x_1 \leq t_1\} \parallel x_2 = t_2$
 $\{x_1 > t_1\} \parallel x_1 = t_3$
3. $\{x_1 > t_3\} \parallel x_2 = t_4$

↳ particija R_1, \dots, R_5 od D

↳ za jednu particiju, $\forall x \in [0,1]^2$,

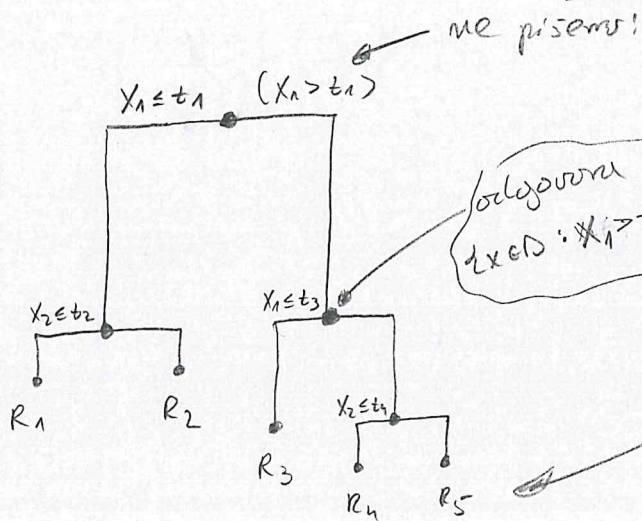
$$\hat{f}(x) := \sum_{m=1}^5 c_m \mathbb{1}_{\{x \in R_m\}} = c_k \text{ ako } x \in R_k,$$

za neke $c_1, \dots, c_5 \in \mathbb{R}$.

• ako $L(y, \hat{y}) = (y - \hat{y})^2 \implies$

$$c_m = \frac{1}{\#\{x^{(i)} \in R_m\}} \cdot \sum_{x^{(i)} \in R_m} y_i$$

↳ ekvivalentan zapis preko binarnog stabla:



odgovora $\{x \in D : x_1 > t_1\} \in D$

[stebila možemo crtati i za $p > 3$!]

$R_1, \dots, R_5 = \text{listovi}$

↳ Prednosti:

- interpretacija [npr. primjene u medicini]
- automatski odabir konjunktata i interakcije među njima
- invarijantnost na monotone transformacije konjunktata

Pluse:

- prediktivna sposobnost → bagging, slučajne šume, boosting

[Konstrukcija stabala u CART algoritmu?]

2.1.1 Regresijska stabla

• $\forall y \in \mathbb{R}$, te postp. $X \in D \subseteq \mathbb{R}^p$ (sva x_j kont.)

Za svako binarno stablo T , svaki čvor $t \in T$ identifikiramo s odgovarajućim područjem od D , te

$$\tilde{T} := \{t \in T : t \text{ je list}\}$$

$$\tau(t) := \{(x^{(i)}, y_i) : x^{(i)} \in t\} \quad (8.1)$$

$$n(t) := |\tau(t)|$$

$$\bar{y}(t) := \frac{1}{n(t)} \sum_{x^{(i)} \in t} y_i, \quad (8.2)$$

$$\hat{A}_T(x) := \sum_{t \in \tilde{T}} \bar{y}(t) \mathbb{1}_{\{x \in t\}}, \quad x \in D$$

[dakle, pretpostavili smo $L(y, \hat{y}) = (y - \hat{y})^2$]

Neka je "cost"

$$C(T) := \frac{1}{n} \sum_{i=1}^n (y_i - \hat{A}_T(x^{(i)}))^2 \quad (= L_T(\hat{A}_T)) \quad (8.3)$$

$$= \frac{1}{n} \sum_{t \in \tilde{T}} \underbrace{\sum_{x^{(i)} \in t} (y_i - \bar{y}(t))^2}_{=: RSS(t)}$$

Kritarij dijeljenja: u svakom listu $t \in \tilde{T}$, neka su

$$t_L = t_L(j, s) := \{x \in t : x_j \leq s\}, \quad j \in \{1, \dots, p\}, s \in \mathbb{R},$$

$$t_R = t_R(j, s) := \{x \in t : x_j > s\}$$

Sve moguće podjele, te $T(j, s)$ rezultirajuće stablo.

↳ tražimo [za s dovoljno glatki samo $n(t)$ unjednakosti!]

$$(j^*, s^*) := \underset{j, s}{\operatorname{arg\,max}} \underbrace{\{C(T) - C(T(j, s))\}}_{=: \Delta C(t, j, s)} \quad (8.4)$$

Nap. 8.2.1

(a) iz (8.3) odmah sledi:

$$\Delta C(t, j, s) = \frac{1}{n} (RSS(t) - RSS(t_L) - RSS(t_R)) \quad (8.5)$$

$$t \in (j^*, s^*) = \underset{j, s}{\operatorname{arg\,min}} \{RSS(t_L) + RSS(t_R)\} \quad (8.6)$$

oni su samo u $T(t)$!

$$\frac{1}{n(t_L)} RSS(t_L) = \widehat{\operatorname{Var}}(y_i : x^{(i)} \in t_L)$$

(b) Vrijedi $\Delta C(t, j, i, s) \geq 0, \forall t, j, i, s$ te (102) 84

$$\Delta C(t, j, i, s) = 0 \Leftrightarrow \bar{y}(t) = \bar{y}(t_L) = \bar{y}(t_0).$$

Kritenij zaustavljanja?

↳ broj listova $|T|$ kontrolira kompleksnost!

↳ N_e dijeli $t \in \tilde{T}$ ako

(1) $\Delta C(t, j, i, s) \leq \delta$ za neki $\delta \geq 0$, [mindeu]
↳ [previše "kretkončno" za velike δ]

ili

(2) $\min \{ n(t_L^*), n(t_0^*) \} < g_1$ [mincut]

za $g_1 \in \mathbb{N}$

ili

(3) $n(t) < g_2$ [minsize]
za $g_2 \in \mathbb{N}$

↳ Kako odrediti δ, g_1, g_2 ?

↳ Alternativa

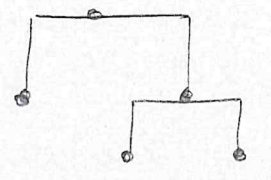
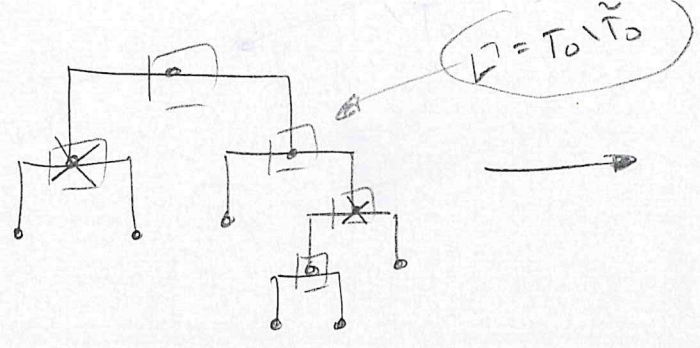
Obrezivanje stabla (engl. "pruning")

• naprimo "veliku" stablo (npr. $\delta=0, g_1=1, g_2=5$ gore)
→ T_0

• $T \subseteq T_0$ je podstablo od T_0 (" $T \subseteq T_0$ ")

ako je dobiveno "obrezivanjem" bilo kojeg broja čvorova $t \in T \setminus T_0$.

mpv.



$T \leq \tilde{T}_0$

T_0

[spajamo neke elemente particije (unatrag)].

• za sve $\alpha \geq 0$,

$C_\alpha(T) := C(T) + \alpha \cdot |\tilde{T}|$ (8.7)

te

$T(\alpha) := \arg\min_{T \leq T_0} C_\alpha(T)$ (8.8)

penalizacija kompleksnosti

Nap. 8.31

(a) $\alpha = 0 \Rightarrow T(0) = T_0$

veći $\alpha \Rightarrow$ manji $|\tilde{T}(\alpha)|$

(b) Može se pokazati da postoji $k \in \mathbb{N}$ te

$0 = \alpha_0 < \alpha_1 < \dots < \alpha_k < \alpha_{k+1} = +\infty$,

i niz stabala $T^0 \geq T^1 \geq \dots \geq T^k$ t.d.

$T(\alpha) = T^i$ za $\alpha \in [\alpha_i, \alpha_{i+1})$, $i=0, \dots, k$ (8.9)

(vidi Ripley (1996))

$(T^0 = T_0, |\tilde{T}^k| = 1)$

(2) $\alpha(t_j, i)$ oclakujemo CV metodom ili na bestnom > kupcu. [86]

[za CV $\rightarrow \forall T_k, k=1, \dots, K$, ponovo računamo

$$T_0 = T_0(T_k), \quad T(x) = T(x, T_k), \quad x \geq 0, \quad (!)$$

te grešku prognoziramo na $T - T_k =$]

California Housing. R

8.1.2 Klasifikacijska stabla

• $\gamma \in \{0, 1, \dots, K-1\} := S$

za stabla T , i sve $t \in T$,

$$m_k(t) := |\{x^{(i)} \in t : \boxed{y_i = k}\}|, \quad (8.10)$$

te

$$\hat{P}(\gamma = k | X \in t) = P_k(t) := \boxed{\frac{m_k(t)}{m(t)}}. \quad (8.11)$$

$$\hat{A}_T(x) := \sum_{t \in \hat{T}} k(t) \mathbb{1}_{\{x \in t\}}, \quad x \in D,$$

te ako je gubitak 0-1, stajemo

$$k(t) = \arg \min_{h \in S} \hat{P}_h(t). \quad (8.12)$$

$$= \arg \min_{h \in S} m_h(t)$$

(8.11)

Kako djelimo čvorove:

kojima, za regresijska stabla,

$$C(T) = \frac{1}{n} \sum_{t \in \tilde{T}} \sum_{x^{(i)} \in t} (y_i - \bar{y}(t))^2$$

$$= \sum_{t \in \tilde{T}} \hat{p}(t) \cdot i(t) \quad (8.13)$$

za

$$\hat{p}(t) := \frac{n(t)}{n} \quad (= \hat{P}(X \in t)) \quad (8.14)$$

$$i(t) := \frac{1}{n(t)} \sum_{x^{(i)} \in t} (y_i - \bar{y}(t))^2$$

$$= \widehat{\text{Var}}(y_i : x^{(i)} \in t)$$

Mjera "nečistoće"
(engl. "impurity")
čvoru t

Za konstrukciju stabla u klasifikaciji koristimo

(1.) Ginijev indeks:

$$i(t) = \sum_{\substack{k \neq k', \\ k, k' \in S}} \hat{p}_k(t) \hat{p}_{k'}(t) = \sum_{k \in S} \hat{p}_k(t) (1 - \hat{p}_k(t)) \quad (8.15)$$

ili:

(2.) Entropija:

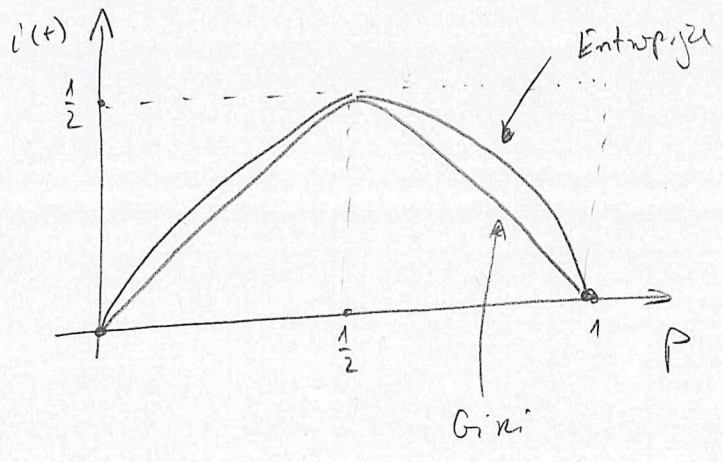
$$i(t) = - \sum_{k \in S} \hat{p}_k(t) \log(\hat{p}_k(t)) \quad (8.16)$$

↳ ujedini • $i(t)$ maksimalna ako $(\hat{p}_0(t), \dots, \hat{p}_{k-1}(t)) = (\frac{1}{k}, \dots, \frac{1}{k})$
 • $i(t)$ minimalna ako $\hat{p}_k(t) = 1$ za neki $k \in S$.

Pr. 8.4 | $S = 20,1\%$ \rightarrow $\left[p := \hat{p}_1(t), p_0(t) = 1-p \right]$ [88]

\Rightarrow Gini $\rightarrow i(t) = 2p(1-p)$

Entropija $\rightarrow i(t) = -p \log(p) - (1-p) \log(1-p)$



Kritenij_djeljanje: za $t \in \tilde{T}$, tražimo

$(j^*, i^*) = \underset{j, i}{\text{argmax}} \{ \Delta C(t, i, j) \}$

(8.17) $= \arg \max_{j, i} \{ \hat{p}(t) i(t) - \hat{p}(t_L) i(t_L) - \hat{p}(t_R) i(t_R) \}$

$= \arg \min_{j, i} \{ \hat{p}(t_L) i(t_L) + \hat{p}(t_R) i(t_R) \}$

Nap. 8.5 | u regresiji imamo $L(y_i, \hat{y}) = (y_i - \hat{y})^2$ te koristimo

$i(t) = \frac{1}{n(t)} \sum_{x^{(i)} \in t} L(y_i, \hat{A}_T(x^{(i)}))$,
 $= \bar{y}(t)$

tj. $C(T) = L_T(\hat{A}_T)$

\rightarrow zašto u klasifikaciji ne koristimo 0-1 gubitak

$i(t) = \frac{1}{n(t)} \sum_{x^{(i)} \in t} \{ \frac{1}{2} y_i + \hat{A}_T(x^{(i)}) \}$
 $= \frac{1}{2} (y_i + \hat{y}(t))$

$$= \boxed{1 - \hat{p}_k(t)} \quad (= \hat{P}(Y \neq \hat{x}_T(x) | X \in t)) \quad (8.10) \quad \boxed{89}$$

mp. $S = 40,13$, $t = (400, 400)$ $\rightarrow n(t) = 800$

\uparrow \uparrow
 $n_0(t)$ $n_1(t)$

podjele (t_L, t_0) :

(10) $t_L = (300, 100)$
 $t_0 = (100, 300)$

(20) $t_L = (200, 400)$
 $t_0 = (200, 0)$

u (10) imamo $\hat{p}(t_L) = \hat{p}(t_0) = \frac{400}{n}$, te

$$\hat{p}(t_L) i(t_L) + \hat{p}(t_0) i(t_0) = \frac{1}{n} \begin{cases} 300, & \text{Gini} \\ 440,9, & \text{Entropija} \\ \boxed{200}, & 0-1 \end{cases}$$

u (20) imamo $\hat{p}(t_L) = \frac{600}{n}$, $\hat{p}(t_0) = \frac{200}{n}$ te

$$\frac{1}{n} \begin{cases} 266,7, & \text{Gini} \\ 381,9, & \text{entropija} \\ \boxed{200}, & 0-1 \end{cases}$$

\Rightarrow 0-1 mjera ne preferira (20) iako daje
 (8.17) "čist" ovor [kosi ne treba dalje dijeliti!]

Ipak, 0-1 mjera konstantna pri obserzaciji stabla. \square

[Zanimanje nas predikcija] 050

8.1.3 Generalizacije

Kategorijalne kovarijate

Ako je $x_j \in \{0, 1, \dots, k_j - 1\} =: S_j$, imamo $2^{k_j - 1} - 1$

možućih podjela S_j na druge grupe

(npr. $t_L = \{x \in t : x_j \in \{0, 1, 3\}\}$, $t_R = \{x \in t : x_j \in \{2, 0, 1, 3\}\}$)

\rightarrow za velike (k_j) imamo pristupaost za izbor x_j za podjela!

Funkcija gubitka \rightarrow pretp. $S = \{0, 1, 3\}$

Umjesto samo $L(\hat{A}) = P(Y \neq \hat{A}(X))$, često nas više zanimaju

• $P(Y = \hat{A}(X) | Y = 1) \rightarrow$ "sensitivity"

• $P(Y = \hat{A}(X) | Y = 0) \rightarrow$ "specificity"

npr. (a) 0 = nije bolest, 1 = bolest

\Rightarrow želimo što veću sens.

(b) 0 = nije spom, 1 = spom

\rightarrow specif. \square

↳ Matrica gubitka je $L = (L_{k,k'} : k, k' \in \{0, 1\})$, gdje je (t-ja)

- $L_{k,k} = 0$, ($\forall k$)

- $L_{k,k'} =$ trošak ako stavimo mjehot (k) predviđeno $\hat{y} = k'$, ($k \neq k'$)

Za stablo T , ako je $x \in t$ ($t \in \tilde{T}$), i

imamo

$$\hat{A}_T(x) = k(t) := \underset{k \in \{0,1\}}{\operatorname{argmin}} \sum_{l \in \{0,1\}} L_{l,k} \hat{P}_l(t) = \hat{P}(Y=l | X \in t)$$

kao u (7.2)

$$= \underset{k \in \{0,1\}}{\operatorname{argmin}} L_{1-k,k} \hat{P}_{1-k}(t) = \hat{E}[L_{Y,k} | X \in t]$$

$L_{k,k} = 0$

to:

$$\hat{A}_T(x) = 1 \iff L_{0,1} \hat{P}_0(t) \leq L_{1,0} \hat{P}_1(t)$$

$$\iff \hat{P}_1(t) \geq \frac{L_{0,1}}{L_{1,0} + L_{0,1}} =: \alpha$$

$\hat{P}_0 = 1 - \hat{P}_1$

↳ • $L_{0,1} > L_{1,0}$ (npr. 0 = nije sporn, 1 = sporn)

$\Rightarrow \alpha > \frac{1}{2}$ te vodi specificity.

• $L_{0,1} < L_{1,0}$ (npr. 0 = nije boleat, 1 = boleat)

$\Rightarrow \alpha < \frac{1}{2}$ te vodi sensitivity

• $L_{0,1} = L_{1,0}$ (0-1 gubitak)

$\Rightarrow \alpha = \frac{1}{2}$.

• uvijek postoji "trade off" između sensit. i specif.

→ ROC krivulja

92

• umjesto kod predikcija, L se može uzeti u obzir i pri konstrukciji stabla