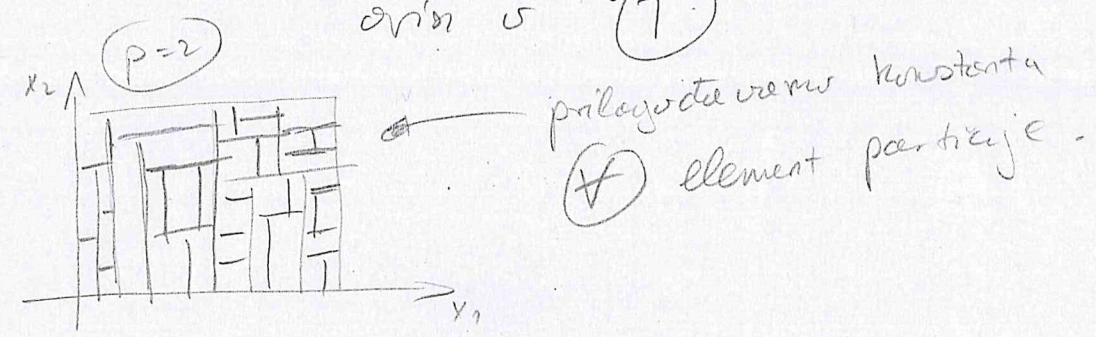


8.2) Bagging i slučajne šume [random forests] → Breiman (1996, 2001)

- Stabla odlučivanja → interpretabilni modeli, ali loša prediktivna sposobnost
- npr., velika stabla → mala prostornost, ali velika varijanca jer je svaki element u (T)



→ ideja je uprosječiti predikciju velikog broja veličitih stabala dobivenih iz bootstrap uzoraka, te tako smanjiti varijancu. Paž, tako direktno gubimo interpretabil.

8.2.1) Bootstrap uzorci [model z_i]

• imamo

$$Z = \{ (x^{(i)}, y_i) : i=1, \dots, n \}$$

koji je realizacija njog uzorka

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

iz $P((X, Y) \in \cdot)$

• za dati (T) , jedan bootstrap uzorak (duljine n), je

$$T_{boot} = \{z_1^*, \dots, z_n^*\},$$

pri čemu su z_1^*, \dots, z_n^* njedi t.d.

$$P(z_1^* = z_i) = \frac{1}{n}, \quad \forall i=1, \dots, n.$$

Čtj., T_{boot} je njedi uzork iz empirijske razdružbe \hat{F}_n uzorku T

(1) Budući da je $\hat{F}_n \approx P((X, Y) \in \cdot)$ za velike n ,

T_{boot} je uzorak koji je približno iz $P((X, Y) \in \cdot)$ [tj. stvarne razdružbe!]

(2.) Neki $z_i \in T$ će se u T_{boot} ponoviti više puta, a neki nit jednom

$\rightarrow T_{boot}$ "drugačiji" od T

mpn! $\forall z_i \in T$,

$$P(z_i \notin T_{boot}) = \underbrace{\left(1 - \frac{1}{n}\right)^n}_{P(z_1^* \neq z_i)} \xrightarrow{n \rightarrow \infty} e^{-1} \approx 0.368$$

te

$$E\left[\frac{1}{n} \sum_{i=1}^n 1_{\{z_i \notin T_{boot}\}}\right] \xrightarrow{n \rightarrow \infty} 0.368 \quad (2.19)$$

\Rightarrow za velike n , T_{boot} u prosjeku izostavlja

36.8% podataka iz T !

[Nap.] Bootstrap metoda se najviše koristi za konstrukciju povrđenih intervala za procjenitelje

$S = S(z_1, \dots, z_n)$ čija je distr. preteška odrediti analitički.

8.2.2 Bagging - (Bootstrap aggregation) 95

1. Napravimo B (B "velik") bootstrap vzorce

$$T_{boot}^{(1)}, \dots, T_{boot}^{(B)}$$

2. $\forall b=1, \dots, B$, konstruiramo (tipično) veliko stabilno $T^{(b)}$ na temelju $T_{boot}^{(b)}$ [bez obzira]

$$\hat{f}_b := \hat{f}_{T^{(b)}}$$

3. Bagging prognetelj je

i) ako $Y \in \mathbb{R}$,

$$\hat{f}_{bag}(x) := \frac{1}{B} \sum_{b=1}^B \hat{f}_b(x), \quad x \in \mathbb{R}^p \quad (8.20)$$

ii) ako $Y \in \{0, 1, \dots, K-1\} =: S$,

$$\hat{f}_{bag}(x) := \operatorname{argmax}_{k \in S} |\{b \in \{1, \dots, B\} : \hat{f}_b(x) = k\}|$$

["majority vote"]

ili

$$\hat{P}_k^{bag}(x) := \frac{1}{B} \sum_{b=1}^B \hat{P}_k^{(b)}(x), \quad x \in \mathbb{R}^p \quad (8.22)$$

te

$$\hat{f}_{bag}(x) := \operatorname{argmax}_{k \in S} \hat{P}_k^{bag}(x), \quad x \in \mathbb{R}^p \quad (8.23)$$

[ako $L = 0-1$]

Prop. 8.4.11z (8.20) i JZVB-a sledi

$$\hat{f}_{bag}(x) \xrightarrow{B \rightarrow \infty} \mathbb{E}_{T_{boot}^{(1)}} [\hat{f}_{T^{(1)}}(x)]$$

[$\hat{f}_{bag}(x)$ je tzv. Monte-Carlo prognetelj!]

(8.24)

\Rightarrow veći B me znači kompleksniji model!

Zašto bagajni - pomaze? [heuristicki]

• prop. $(Y \in \mathbb{R})$, te $\hat{\lambda} := \hat{\lambda}_{T_0}$ za vel. k steper $T_0 = T_0(T)$
 ↑
 mig. uvest.

↳
$$E_T[L_x(\hat{\lambda})] = d^2 + \text{Var}_T(\hat{\lambda}(x)) + \text{Bias}_T[\hat{\lambda}(x)]$$

\approx (prop.) $d^2 + \underbrace{\text{Var}_T(\hat{\lambda}(x))}_{=: d_{T_0}^2}$

• $\hat{\lambda}_A(x), \dots, \hat{\lambda}_B(x)$ su jednako distrib. (dla znanje!),
 te prop. da je jer onke o T!

$$\text{Var}(\hat{\lambda}_b(x)) \approx d_{T_0}^2, \text{Bias}(\hat{\lambda}_b(x)) \approx 0, (8.25)$$

Najbolje je od
 $T_{boot}^{(A)}, \dots, T_{boot}^{(B)}$!

te neka je

$$p(x) := P(\hat{\lambda}_b(x), \hat{\lambda}_{b'}(x)), \quad b \neq b' \quad (8.26)$$

↳ [konkluzija]

↳ (1)
$$E[\hat{\lambda}_{avg}(x)] = \frac{1}{B} \sum_{b=1}^B E[\hat{\lambda}_b(x)]$$

$= E[\hat{\lambda}_1(x)],$ (8.25)

↳
$$\text{Bias}(\hat{\lambda}_{avg}(x)) = \text{Bias}(\hat{\lambda}_1(x)) \approx 0.$$

(2)
$$\text{Var}(\hat{\lambda}_{avg}(x)) \stackrel{(8.25)}{=} \left(p(x) + \frac{1-p(x)}{B} \right) \cdot \text{Var}(\hat{\lambda}_1(x))$$

\approx ——— || ——— $d_{T_0}^2$

(8.25)

2
zu
velike B

$$p(x) \cdot G_{T_0}^2$$

$$[\Rightarrow p(x) \approx 0] \quad (8.27)$$

$$G_{T_0}^2$$

$$p(x) < 1$$

8.2.3 Slučajne šume (random forests) \rightarrow (RF)

ideja je manjiti $p(x)$ u (8.25)

u odnosu na bagging, jedina promjena je u konstrukciji stabla $T^{(b)}$ iz $T_{\text{root}}^{(b)}$, $\forall b=1, \dots, B$:

Prje svakoj dijeljenju čvora, izaberu slučajno (m) korijata ($m \leq p$) kao kandidate za dijeljenje.

intuitivno, manji $(m) \Rightarrow$ manja $p(x)$

[stabla primorana koristiti različite korijate za dijeljenje]

(ipak, mjeruju se i $\text{Var}(\hat{f}_B(x))$; $\text{Bias}(\hat{f}_B(x))$)

\rightarrow [vidi ESL, Page. 15.4]

(m) je najvažniji hiperparametar. ("mtry")

[RF često rade odlično, a jednostavno su za prilagodbu ("off-the-shelf")]

Out-of-Bag ("OOB") procjena greške

- $\forall b \in \{1, \dots, B\}$,

$$\mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{z_i \notin T_{boot}^{(b)}\}} \right] \stackrel{(8.14)}{\approx} 0.368$$

"OOB uvijek" \Rightarrow misu koristiti za konstr. \hat{f}_b !

- $\forall i=1, \dots, n$

$$OOB_i := \{ b \in \{1, \dots, B\} : z_i \notin T_{boot}^{(b)} \}, \quad (8.28)$$

- ako je $(y \in \mathbb{R})$, defin.

$$\hat{y}_i^{OOB} := \frac{1}{|OOB_i|} \sum_{b \in OOB_i} \hat{f}_b(x^{(i)}), \quad i=1, \dots, n$$

te OOB procjenu greške

$$OOB(\hat{f}_{\text{fit}}) := \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i^{OOB})^2 \quad (8.29)$$

\hookrightarrow slične procjene kao i LOOCV, ali dobivamo ih besplatno, tj. u sklopu treniranja!

- ako $y \in \{0, \dots, K-1\} =: S$,

$$\hat{y}_i^{OOB} := \operatorname{argmax}_{k \in S} |\{ b \in OOB_i : \hat{f}_b(x^{(i)}) = k \}|,$$

$$te \quad OOB(\hat{f}_{\text{fit}}) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{y_i \neq \hat{y}_i^{OOB}\}}. \quad (8.30)$$

Uticaj kovarijete (engl. variable importance) ($\forall \in \mathbb{R}$) 95

[• \hat{A}_{ort} (i specijalno, \hat{A}_{maj}) više nemaju strukturu stabla pa se gubi njihova sposobnost interpretacije!]

• Za jedno stablo T , neka je $(\forall j = 1, \dots, |P|)$,

"importance"

$$I_j(T) := \sum_{\substack{t \in T \setminus \tilde{T}: \\ j^*(t) = j}} \Delta C(t, j^*(t), s^*(t)) \quad (8.31)$$

uz

• $X_{j^*(t)}$ je kovarijeta korištena za djeljenje čvora t ,

• a $\Delta C(t, j^*(t), s^*(t))$ je omjer "čvorka" C (vidi (8.4)) pri djeljenju čvora (t) .

• Za RF definiramo

$$I_j := \frac{1}{B} \sum_{b=1}^B I_j(T^{(b)}) \quad (8.32)$$

Gratni djelomične zavisnosti (partial dependence plots) 100

• pref. $Y \in \mathbb{R}$

↳ (PDP)

• za $A, B \subseteq \{1, \dots, p\}$ t.d. $A \cap B = \emptyset$, $A \cup B = \{1, \dots, p\}$, $\forall x \in \mathbb{R}^p$

pišemo $x_A := (x_j)_{j \in A}$, $x_B := (x_j)_{j \in B}$

• za $f := f_{\text{RF}} : \mathbb{R}^p \rightarrow \mathbb{R}$, projektnu ili djelomičnu zavisnost f -je f na X_A defin. kao

$$f_A(x_A) := \mathbb{E} [f(x_A, X_B)], \quad x_A \in \mathbb{R}^{|A|} \quad (8.33)$$

↑
briše, $|A|=1$ & 2

↑
dućajms

[f_A predstavlja utjecaj kovarijata X_A na $f(x)$ independently isto smislu u obzir uzeti projekcion utjecaj ostalih kovarijata X_B na $f(x)$.]

(mfr.) - ako je $f(x) = h_1(x_A) + h_2(x_B)$,

$$f_A(x_A) = \boxed{h_1(x_A)} + \underbrace{\mathbb{E} [h_2(X_B)]}_{= \text{const. !}}$$

- ako je $f(x) = h_1(x_A) \cdot h_2(x_B)$,

$$f_A(x_A) = \text{const.} \cdot h_1(x_A)$$

• $f_A(x_A)$ projekcijeno > (8.34)

$$\hat{f}_A(x_A) := \frac{1}{m} \sum_{i=1}^m f(x_A, x_B^{(i)}) \quad (8.34)$$

[za RF moguće efikasno računati] [(8.33) nije ograničen samo na RF!]

Nap. | Ako su X_A i X_B mezanisne,

100

$$E[f(X_A, X_B) | X_A]$$

$$= f_A(X_A), \text{ tj.}$$

$$f_A(x_A) = E[f(\underbrace{X_A, X_B}_{=X}) | X_A = x_A]$$

→ Boston - bag - i - rt. R

općenito, to ne možemo
procijeniti jer nemamo
dovoljno $(x^{(i)}, y_i) \in \mathcal{T}$
t.d. $x_A^{(i)} = x_A$