

### 8.3 Boosting -

↳ ideja: kombinirati sekvencijalno niz jednostavnih "baza" procjenitelja u jedan (jako) dobar procjen., pri čemu se svaki novi bazni procjenitelj više fokusira na elemente iz  $I$  koji su lošije procijenjeni u trenutnom procjenitelju.

- bazni procjenitelji mogu biti bilo što
- mi ćemo promatrati Friedmanov (2001) "Gradient Boosting Machine" (GBM) algoritam uz regresijska stabla kao bazne procjenitelje
- metode poput XGBoost i LightGBM su varijante originalnog GBM algoritma, a predstavljaju jedne od najmoćnijih prediktivnih metoda u strojnom učenju

### 8.3.1

### Gradient boosting - za regresiju ( $Y \in \mathbb{R}$ ) ("GB")

• za danu  $\ell$ -ju gubitka  $L: \mathbb{R}^2 \rightarrow [0, \infty)$ , želimo procijeniti  $f^*(x) := \underset{c \in \mathbb{R}}{\operatorname{argmin}} \mathbb{E}[L(Y, c) \mid X=x], x \in \mathbb{R}^p, (8.3)$

na temelju  $I = \{(x^{(i)}, y_i) : i=1, \dots, n\}$ .  
 $=: L_I(x)$

↳ tipično

$$\hat{f} := \underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, f(x^{(i)})) \quad (8.36)$$

neka funkcija

pr. čemu (8.36) rješavamo nekom optimizačjskom metodom,  
 a time da  $L$  u (8.35) i (8.36) ne mora nužno biti ista!

153

Pr. 8.71

(a)  $L(y, f(x)) = \frac{1}{2}(y - f(x))^2$  ( $L_2$ -gubitak) (8.37)

$\Rightarrow f^*(x) = E[Y | X=x]$

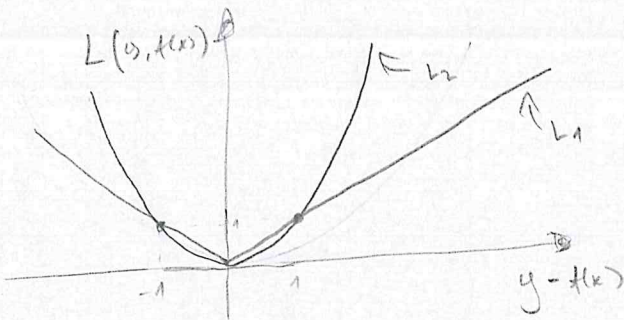
$\Rightarrow$  pogodna za temeljnu analizu

$\Rightarrow$  ipak, na  $L_E(t)$  u (8.36) značajno utječu  $(x^{(i)}, y_i)$  t.d. je  $|y_i - f(x^{(i)})|$  velik ("outlier")

(b)  $L(y, f(x)) = |y - f(x)|$  ( $L_1$ -gub.) (8.38)

$\Rightarrow f^*(x) = \text{median}(Y | X=x)$

$\Rightarrow$  "robustnija" alternativa za  $L_2$ -gub.



$\Rightarrow$  ipak, ako  $y - f(x) \gg 0$ ,  
 $L_1(y, f(x)) \gg L_2(y, f(x))$ ,

(c)  $L_\delta(y, f(x)) = \begin{cases} \frac{1}{2}(y - f(x))^2, & \text{ako } |y - f(x)| \leq \delta, \\ \delta(|y - f(x)| - \frac{1}{2}\delta), & \text{inače} \end{cases}$  (8.39)

za  $\delta > 0$ . ("Huberov" gub.)

$\Rightarrow$  kombinacija  $L_1$  i  $L_2$  gub.: za male  $y - f(x)$  se ponaša kao  $L_2$ , a za velike kao  $L_1$  gub.

[izbor  $L$ -a u konkretnom problemu i još to hoćemo više detalizirati!]

• GB metoda (za stabla):

(i)  $f^*$  proširujemo modelom oblika

$$f^{(M)}(x) = f^{(0)}(x) + \sum_{m=1}^M b_m(x), \quad x \in \mathbb{R}^P \quad (8.40)$$

[tj.  $f$  u (8.36) je familija  $f$ -ju oblika (8.40)]

gdje su  $b_1, \dots, b_M$  regresijska stabla iz (8.1.1), tj.

$\forall m, \exists$  stablo  $T_m$  t.d.

•  $|\tilde{T}_m| = \mathcal{J}$  (tipično  $\mathcal{J}$  mali!), te

•  $b_m(x) := \sum_{t \in \tilde{T}_m} \gamma_t \mathbb{1}_{\{x \in t\}}, \quad \forall x$

$$\gamma_t = \underset{\gamma \in \mathbb{R}}{\operatorname{argmin}} \sum_{x^{(i)} \in t} L(y_i, \gamma) \quad (8.41)$$

(ovdje  $L!$ )

m.p.r. | (a)  $L = L_2 \Rightarrow \gamma_t = \bar{y}(t)$  (kao u (8.1.1))

(b)  $L = L_1 \Rightarrow \gamma_t = \operatorname{median}(y_i : x^{(i)} \in t)$

(ii) rješenje problema (8.36) aproksimiramo

rekurentno po uzoru na metodu

gradijentnog spusta [eng. gradient descent]:

$\forall m = 1, 2, \dots, M$ ,  $b_m$  biramo tako da

$$b_m \approx \underset{b}{\operatorname{argmin}} L_T(f^{(m-1)} + b) \quad (8.42)$$

$$= \left[ - \text{||} - \sum_{i=1}^n L(y_i, f^{(m-1)}(x^{(i)}) + b(x^{(i)})) \right]$$

# Gradientni spust [digracija]

• Neka je  $L: \mathbb{R}^n \rightarrow \mathbb{R}$  diferencijabilna i konveksna

$f$ -ice, a tražimo

$$\Theta^* := \underset{\Theta \in \mathbb{R}^n}{\text{argmin}} L(\Theta)$$

• Algoritam:

1. odredi  $\Theta^{(0)}$  [npr.  $\Theta^{(0)} := 0$ ]

2.  $\forall m=1, \dots, H,$

$$\Theta^{(m)} = \Theta^{(m-1)} + \eta_m \cdot g_m$$

gdje je

$$(a) g_m := -\nabla L(\Theta^{(m-1)})$$

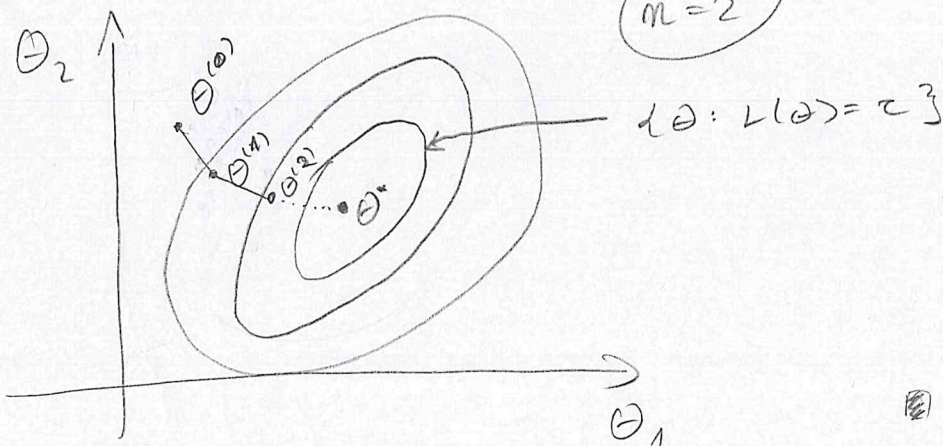
Dugor najvišeg pada (koliko)  $f$ -je  $L$  u  $\Theta^{(m-1)}$   
(negativni gradijent)

$$= \left( \frac{-\partial L}{\partial \theta_1}(\Theta^{(m-1)}), \dots, \frac{-\partial L}{\partial \theta_n}(\Theta^{(m-1)}) \right) \in \mathbb{R}^n$$

$$(b) \eta_m := \underset{\eta > 0}{\text{argmin}} L(\Theta^{(m-1)} + \eta \cdot g_m)$$

3. Vrați  $\Theta^{(H)}$

("line search")



• Shvatišuv sader ovake  $t$ -ju  $f: \mathbb{R}^p \rightarrow \mathbb{R}$  kao vektor

$f \in \mathbb{R}^n$  uz predikcijuje

106

$$f \mapsto (f(x^{(1)}), \dots, f(x^{(n)})) \\ =: (t_1, \dots, t_n),$$

a  $L_{\mathbb{Z}}$  kao  $f$ -ju na  $\mathbb{R}^n$  uz

$$(8.43) \quad L_{\mathbb{Z}}(f) := \sum_{i=1}^m L(y_i, t_i), \quad f = (t_1, \dots, t_n) \in \mathbb{R}^n$$

( $\rightarrow$  ključna ideja: u (8.42) bismo bismo trebali da  
 $(b_m(x^{(1)}), \dots, b_m(x^{(n)}))$  reprezentiramo

$$-\nabla L_{\mathbb{Z}}(f^{(m-1)}) = (-g_{1,m}, \dots, -g_{n,m}) \\ = g_m$$

gdje je  $\forall i=1, \dots, n,$

$$g_{i,m} := \frac{\partial L_{\mathbb{Z}}}{\partial t_i}(f) \Big|_{f=f^{(m-1)}}$$

$$(8.43) \quad \frac{\partial L(y_i, t_i)}{\partial t_i} \Big|_{t_i = f^{(m-1)}(x^{(i)})},$$

pri čemu stablo  $T_m$  konstruiramo koristeći

CART algoritam iz (8.1.1) na početku

$$d(x^{(i)}, [-g_{i,m}]) : i=1, \dots, n \}$$

(dakle, uz  $L=L_2$ )

$\exists$  jako  
 efikasne  
 implem.!!

• "Gradient tree boosting" algoritam

1. Inicijaliziraj  $f^{(0)}$ , [mp.  $\equiv$  const.]

Friedman (2001)

2.  $M = 1, \dots, M$ ,

157

(a)  $\forall i = 1, \dots, n$ , izračunaj

$$g_{i,m} := \frac{\partial L(y_i, f_i)}{\partial f_i} \Big|_{f_i = f^{(m-1)}(x^{(i)})}$$

(b) Konstruiraj regresijski stabler  $\tilde{T}_m$  (kao u (8.1.1)) za podatke

$$\{(x^{(i)}, \underbrace{-g_{i,m}}_{\text{"pseudo-residuals"}}) : i = 1, \dots, n\}$$

t.d.  $|\tilde{T}_m| = \underbrace{J}$ ,

"pseudo-residuals"

(c)  $\forall t \in \tilde{T}_m$ , izračunaj

$$\gamma_t := \underset{\gamma \in \mathbb{R}}{\operatorname{argmin}} \sum_{x^{(i)} \in t} L(y_i, \gamma),$$

"line search"

te defin.

$$b_m(x) := \sum_{t \in \tilde{T}_m} \gamma_t \mathbb{1}_{\{x \in t\}}, \quad x \in \mathbb{R}^p$$

(d)  $f^{(m)}(x) := f^{(m-1)}(x) + \underbrace{\lambda}_{\text{learning rate}} \cdot b_m(x), \quad x \in \mathbb{R}^p$

3. Vraći  $f^{(M)}$ .

gbm paket u R-u

$T$  ni su hiperparametri:

(i)  $(M)$  — broj stabala. Za velike  $M$  moguć overfitting,

ali u praksi se nijetko događa.

CV metode za odabir

[još nije stvar jasen ovaj fenomen!]

ii)  $\lambda \in (0, 1]$  - "shrinkage" parameter [regularizacija!]

• manji  $\lambda \Rightarrow$  veći "optimalni"  $M$ ,

108

• te tipično i bolja prediktivna sposobnost!

[aly. "opornje" uči]

$\rightarrow$  problem ako je  $M$  prevelik

[• tipično binarno  $\lambda$  što manji, npr. iz  $(0.001, 0.1)$ ]

iii)  $J =$  broj listova u svakom stablu.

• kontrolna intenzivacija

•  $\forall T_m$  imamo ukupno  $(J-1)$  okupljanja

$\Rightarrow$  moguće interakcije od najviše

$(J-1)$  kovarijata

• tipično  $J \in \{4, \dots, 8\}$

Def. 8.8 | Po uzoru na slučajne šume, u svakom koraku

$m = 1, \dots, M$ , korake (a) - (c) možemo provesti na slučajnoj

odabranom podskupu od  $T$  veličine  $(p \cdot m)$

za  $p \in (0, 1]$  ("Stochastic GB")

[• default je  $(p = \frac{1}{2})$ ]

Pr. 8.05

$$(a) L(y, t(x)) = \frac{1}{2} (y - t(x))^2$$

$$\Rightarrow -g_{i,m} = \frac{\partial L(y_i, t)}{\partial t} \Big|_{t=t^{(m-1)}(x^{(i)})} = y_i - t^{(m-1)}(x^{(i)})$$

"Obični" residuali za  $t^{(m-1)}$

$\rightarrow$  "L2-boosting"

$$(b) L(y, f(x)) = |y - f(x)|$$

$$\Rightarrow -g_{i,m} = \text{sgn}(y_i - t^{(m-1)}(x^{(i)}))$$

$$= \begin{cases} +1 & , y_i > t^{(m-1)}(x^{(i)}) \\ -1 & , < \end{cases}$$

↑ [uopće ne ovisi o  $|y_i - t^{(m-1)}(x^{(i)})|$   
("robustnost")]

$$(c) L(y, f(x)) = \text{Huber iz (9.34)}$$

$$\Rightarrow -g_{i,m} = \begin{cases} y_i - t^{(m-1)}(x^{(i)}), & \text{ako } |y_i - t^{(m-1)}(x^{(i)})| < \delta \\ \delta \cdot \text{sgn}(y_i - t^{(m-1)}(x^{(i)})), & \text{inače} \end{cases}$$

↳  $\delta$  se tipično bira iz  $\mathbb{Z}$  posredno  $\neq m$ .

### Nap. 8.10 Interpretabilnost?

↳ utjecaj konjugata te glatko parcijalne  
zavisnosti redimov isto kao za RF.

↳  $l_2$ -boosting - primjer

↳ Boosting - boosting.R



8.3.2 GIB za binarnu klasifikaciju

$Y \in \{0, 1\} =: S$

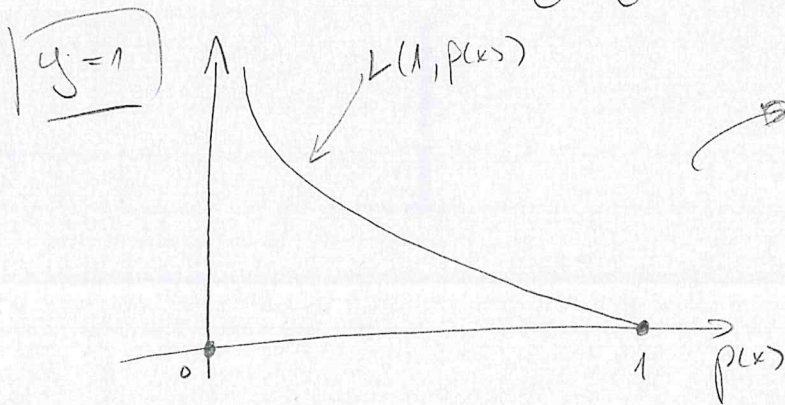
želimo procijeniti:  $p^*(x) := P(Y=1 | X=x), \forall x \in \mathbb{R}^p$

tzv. Bernoullijer gubitak je  $L: S \times [0, 1] \rightarrow [0, \infty)$   
 je  $\rightarrow$  negativna log-vjerojat.

$$L(y, p(x)) = -\log(P(B(p(x)) = y))$$

$$= \begin{cases} -\log p(x), & \text{akr } y=1 \\ -\log(1-p(x)), & \text{akr } y=0 \end{cases}$$

$$= -y \log p(x) - (1-y) \log(1-p(x)) \quad (8.43)$$



$L(1, p(x)) > 0$  i  
 u slučaju  $p(x) \geq \frac{1}{2}$ ,  
 tj. kada je  
 $\hat{f}(x) := \mathbb{1}_{\{p(x) \geq \frac{1}{2}\}} = 1 = y!$   
 $\rightarrow$  želimo  $p(x)$  što bliže  $\textcircled{1}$ !

za to pokušajte da je

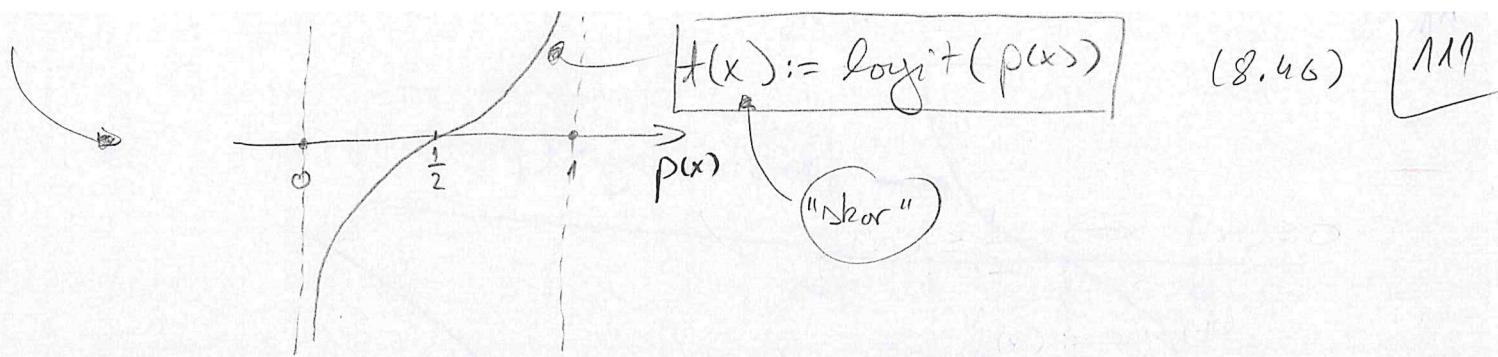
$$p^*(x) = \underset{z \in [0, 1]}{\text{argmin}} \mathbb{E}[L(Y, z) | X=x] \quad (8.44)$$

kar u log-regresiji!

ekvivalentno, umjesto  $p^*(x)$  možemo modelirati

$$f^*(x) := \log \frac{p^*(x)}{1-p^*(x)} \in \mathbb{R} \quad (8.45)$$

(= logit( $p^*(x)$ ))



Takoder,

$$p(x) = \frac{e^{t(x)}}{1 + e^{t(x)}} = \frac{1}{1 + e^{-t(x)}} \quad (8.47)$$

$$1 - p(x) = \frac{1}{1 + e^{t(x)}} \quad (8.48)$$

Bernoullijer yubitatek (o optimom na t) je

$$L(y, t(x)) = -y \log p(x) - (1-y) \log(1-p(x))$$

$$\stackrel{(8.47)}{=} \stackrel{(8.48)}{=} \left[ -y t(x) + \log(1 + e^{t(x)}) \right] \quad (8.49)$$

a iz (8.49) sledi da je

$$t^*(x) = \log \frac{p(x)}{1-p(x)} \stackrel{(\circledast)}{=} \underset{\tau \in \mathbb{R}}{\operatorname{argmin}} \mathbb{E}[L(y, \tau) | X=x] \quad (8.50)$$

$=: t_{\text{Bern}}^*(x)$

$\Rightarrow$  ako u GIB algoritmu iz (8.3.1) vzamemo  $L$  iz (8.49), projekcijem  $t^*(x) = \log \frac{p(x)}{1-p(x)}$ , a pseudo-residuali su

$$-g_{i,m} = \left. \frac{-\partial L(y_i, t)}{\partial t} \right|_{t=t^{(m-1)}(x^{(i)})}$$

$$\stackrel{(8.47)}{=} \stackrel{(8.49)}{=} \left[ y_i - p^{(m-1)}(x^{(i)}) \right] \quad (8.51)$$

• ako je  $f$  progama za  $f^*(x)$ , defn.  $\hat{f}: \mathbb{R}^P \rightarrow S$  da [112]

$$\hat{f}(x) := \mathbb{1}_{\{f(x) \geq 0\}} \quad (= \mathbb{1}_{\{p(x) \geq \frac{1}{2}\}}) \quad (8.52)$$

• Za Bern. gubitak vrijedi  
(8.47)  $\rightarrow$  (8.48)

$$L(y, f(x)) \stackrel{\substack{\text{log} \\ \in \{0,1\}}}{=} (y) \log(1 + e^{-f(x)}) + (1-y) \log(1 + e^{f(x)}) \\ = \log(1 + e^{-\tilde{y} \cdot f(x)}), \quad (8.53)$$

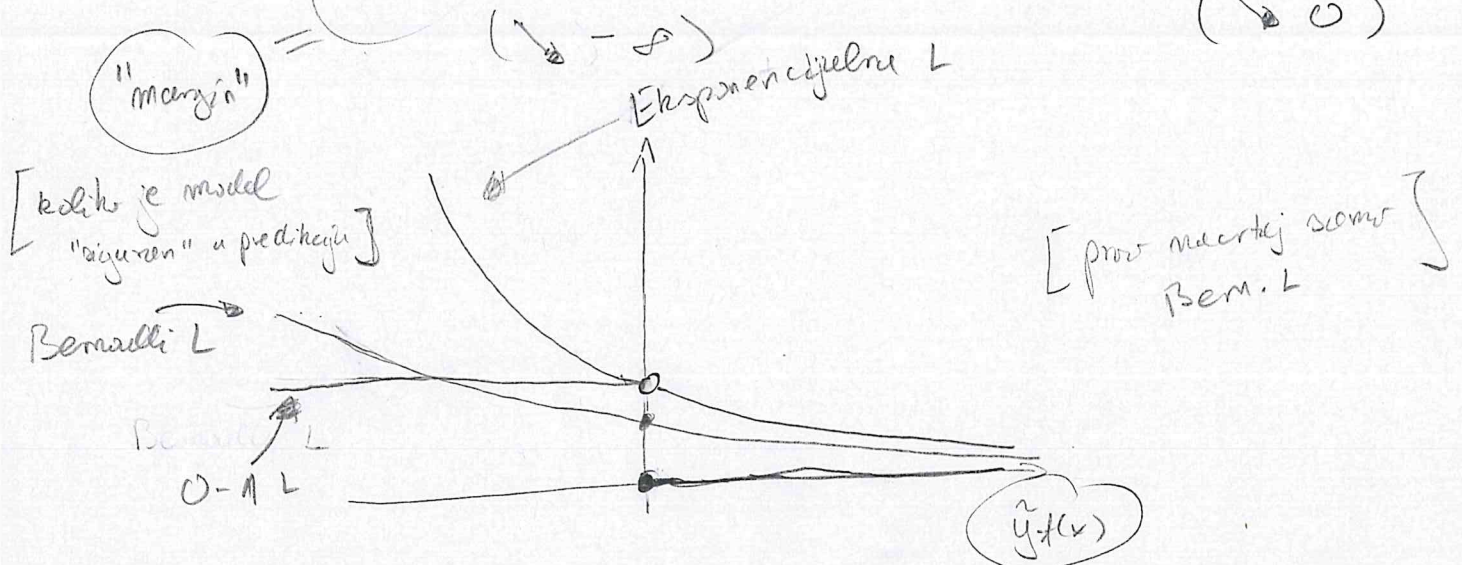
gdje je  $\tilde{y} := 2y - 1 \in [-1, 1]$ .

$\hookrightarrow$  uočimo,

$$f(x) = y \stackrel{(8.52)}{\Leftrightarrow} \tilde{y} \cdot f(x) \geq 0$$

te

$$\tilde{y} \cdot f(x) \geq 0 \Rightarrow p_{\tilde{y}}(x) = \begin{cases} p(x), & y=1 \\ 1-p(x), & y=0 \end{cases} \begin{matrix} \nearrow 1 \\ \searrow 0 \end{matrix}$$



• alternativna za Bern. gubitak je  $f$ -je  $L$  t.d.  $\forall f: \mathbb{R}^P \rightarrow \mathbb{R}$ ,

$$L(y, f(x)) = e^{-\tilde{y} \cdot f(x)} \quad (8.54)$$

"Eksponencijalni gubitak"

(i) iz (8.53) imamo da za Bern. gubitak vrijedi:

$$\tilde{y}_t(x) \rightarrow -\infty \Rightarrow L(y, t(x)) \sim \frac{-\tilde{y}_t(x)}{\text{linearni!}}$$

[primjeri na kraju t jako gubi!]

[malo primjeri na kraju]

Eksponencijalni gub. ima veći problem s "outlierima" [npr. ako imamo greške u podacima] [ili kao L2-gub. za regresiju]

S druge strane, iz (8.53) imamo

$$\tilde{y}_t(x) \rightarrow +\infty \Rightarrow L(y, t(x)) \sim e^{-\tilde{y}_t(x)}$$

(Bern.)  
log(1+x) ~ x, kada x -> 0

(isto kao Eksp. L)

[malo eksp. L na slici].

(ii) Oba gubitka su monotna, konvexna i diferencijabilna aproksimacije 0-1 gubitka

$$L(y, t(x)) = \mathbb{1}_{\{\tilde{y}_t(x) < 0\}}$$

Co [dodaj 0-1 gub. na slici]

Nap. 8.12 GIB uz ekponencijalni gubitak je ekvivalentan

tzv. Adaboost algoritmu (Freund, Schapire, 1996)

→ prvi uspješni boosting algoritmi

→ iako je malo robustan od GIB algoritma uz Bern. gubitak, postaje jako efikasne

implementacije → vidi npr.

Viola-Jones algoritam za detekciju lica

(vidi ES, prof. 10)

→ to da Adekost zuprav optimizira eksperimentalni gubitak je u originalnom algoritmu bilo skom 114  
implicitno → Friedman et. al. (2000)

to pokazuje, te je to zuprav dovelo do GB algoritma u Friedman (2001)! 1

• za DE pokazite da za eksp. gubitak vrijedi:

$$x_{\text{Ehsp}}^* := \underset{c \in \mathbb{R}}{\text{argmin}} \mathbb{E}[L(y, c) | x=x] = \left(\frac{1}{2}\right) \sigma^2(x) \quad (8.55)$$

⇒ ako je  $\sigma$  progrena za  $x_{\text{Ehsp}}^*$ ,  $(2\sigma)$  je

progrena za  $\sigma$ , tj.  $\frac{1}{1 + e^{-2\sigma(x)}}$  je

progrena za  $p^*(x)$  (vidi (8.47)).

↳ specijalno, opet stavljamo

$$f(x) := \mathbb{1}_{\{2\sigma(x) \geq 0\}} = \left[ \frac{1}{2} \sigma(x) \geq 0 \right]$$

⇒ GB algoritam uz eksp. gubitak progrena je

$x_{\text{Ehsp}}^*$ , a pseudo-residuali su

$$-g_{i,m} \stackrel{\text{DE}}{\approx} \tilde{y}_i e^{-\tilde{y}_i \sigma^{(m-1)}(x^{(i)})} \quad (8.56)$$

$$\hookrightarrow \tilde{y}_i \sigma^{(m-1)}(x^{(i)}) < 0 \Rightarrow | -g_{i,m} | = e^{|\tilde{y}_i \sigma^{(m-1)}(x^{(i)})|}$$

↳ [prije klasifikirani primjeri]

[koliko smo bili "sigurni" u klasifikaciju]

Rep.

- značajnost vanjske i unutari parcijalne odvoda (za  $f^{(n)}$ , me za  $\pi^{(n)}$ ) kao i za regresiju

- u slučaju  $Y \in \{0, \dots, k-1\}$ , za  $\underline{k \geq 3}$ , modeliramo

(k)  $f$ -je  $(t_0, \dots, t_{k-1})$  kao u multinomijalnoj log.

regresiji: npr.  $t_0 \equiv 0$ , a za  $k \geq 1$

$$f_k(x) = \log \frac{P(Y=k | X=x)}{P(Y=0 | X=x)}$$

$$P(Y=k | X=x) = \frac{e^{t_k(x)}}{\sum_{k=0}^{k-1} e^{t_k(x)}}$$

[izostavljam detalje]

- ako je  $Y \in \{0, 1, 2, \dots\} = \mathbb{N}_0$ , kao u GLM-u, za Poissonom raspodjelu, GLS algoritmom možemo

modelirati:

$$f(x) := \log \{E[Y | X=x]\} \quad (\in \mathbb{R}!)$$

te koristiti:

$$L = - \log\text{-vjerovatnost}$$