

Lecture 8



• CART: Classification & Regression Trees

- Recursive partitioning $(R_1, \dots, R_m, \dots, R_M)$



Loss $(R_1, \dots, R_m, \dots, R_M)$
purity:

• Classification:

$\{1, \dots, c, \dots, C\}$ categories

in R_m , $n_m = \#$ examples in R_m
 $n_{m,c} = \#$ " " " in c

$$\hat{p}_{m,c} = \frac{n_{m,c}}{n_m}$$

$$\begin{aligned} \text{Loss}(R_1, \dots, R_m, \dots, R_M) &= - \sum_{m=1}^M n_m \sum_{c=1}^C \hat{p}_{m,c} \log \hat{p}_{m,c} \\ &= \sum_{m=1}^M n_m \text{entropy}(\hat{p}_m) \end{aligned}$$

	1	2	...	c	...	C
\hat{p}_m	$\hat{p}_{m,1}$	$\hat{p}_{m,2}$...	$\hat{p}_{m,c}$...	$\hat{p}_{m,C}$

empirical distribution in R_m

- Regression:

$$\text{Loss}(R_1, \dots, R_m, \dots, R_M) \\ = \sum_{m=1}^M \sum_{i: x_i \in R_m} (y_i - \hat{c}_m)^2$$

$$\hat{c}_m = \text{average within } R_m$$

$$= \frac{\sum_{i: x_i \in R_m} y_i}{n_m}$$

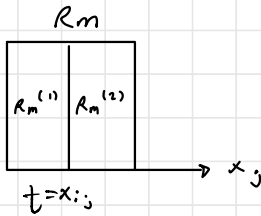
- Recursive Partition:

In step M $\{R_1, \dots, R_m, R_M\}$
for m in $1:M$

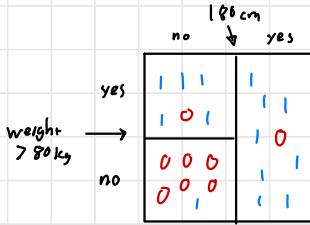
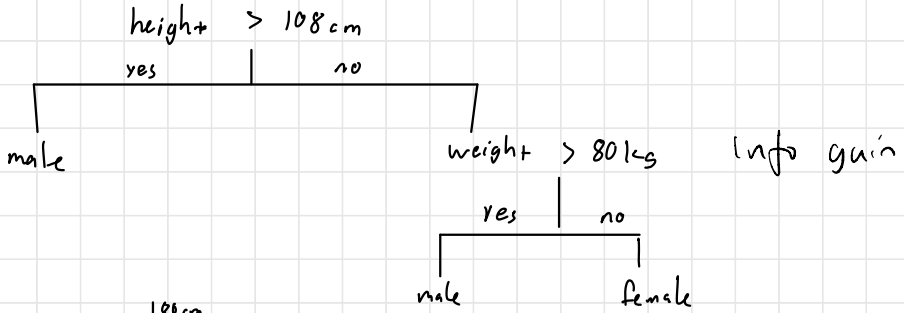
for j in $1:P$

for t in $\{x_{i_j}, i=1, \dots, n\}$

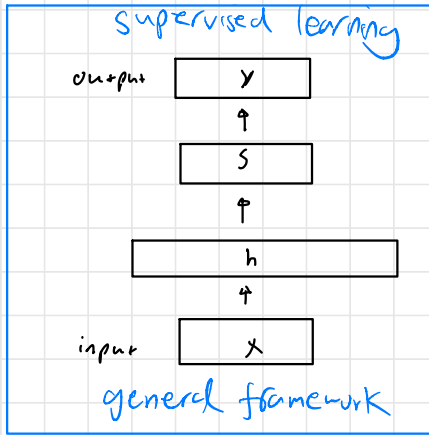
info-gain = $\text{Loss}(R_1, \dots, R_m, \dots, R_M) - \text{Loss}(R_1, \dots, R_m^{(1)}, R_m^{(2)}, \dots, R_M)$



o Example:



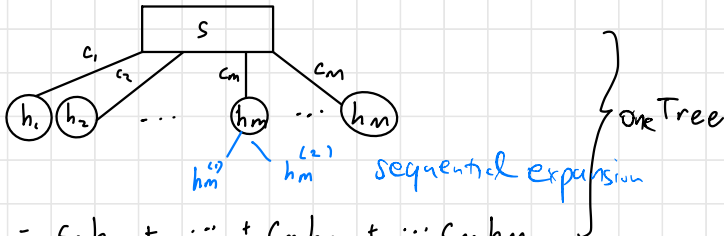
o Recall:



Regression: $y \sim \mathcal{N}(s, \sigma^2 I)$

Classification: $y \sim \text{softmax}(s)$

$$P(y=k|s) = \frac{e^{s_k}}{\sum_c e^{s_c}}$$



$$s = c_1 h_1 + \dots + c_m h_m + \dots + c_n h_n$$

$$h_m = \mathbb{1}(x \in R_m)$$

- We may express in matrix form:

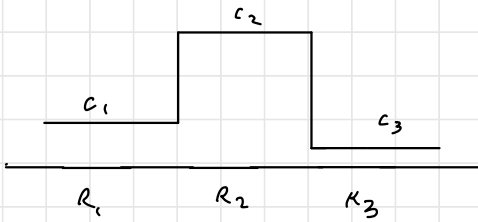
$$s = \begin{pmatrix} c_1 & \dots & c_m & \dots & c_M \end{pmatrix} \begin{pmatrix} h_1 \\ \vdots \\ h_m \\ \vdots \\ h_M \end{pmatrix}$$

$$(= w h)$$

↑
one-hot vector

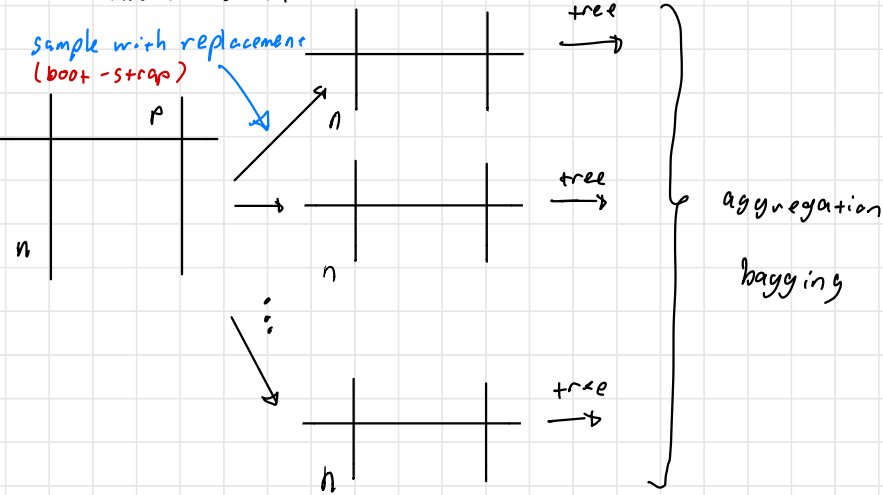
- If $h_m \geq 1$, i.e. $x \in R_m$ then $s = c_m$

$$s_0 s = f(x) = c_m \text{ if } x \in R_m$$

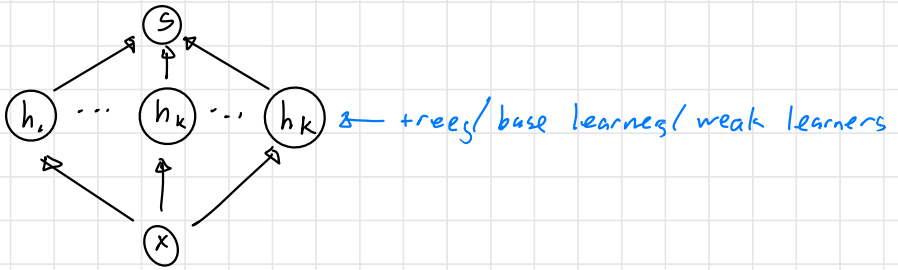


• Random Forest:

sample with replacement
(boot-strap)



• Boosting



$$S = f(x) = \sum_{k=1}^K h_k(x)$$

• Extreme - Gradient - Boosting :

$$\text{Log-likelihood} = \sum_{i=1}^n \ell(s_i) = \sum_{i=1}^n \ell(\hat{s}_i) + \ell'(\hat{s}_i) \Delta s_i + \frac{1}{2} \ell''(\hat{s}_i) \Delta s_i^2$$

$$s_i = \underbrace{\sum_{k=1}^{k-1} h_k(x_i)}_{\hat{s}_i} + \underbrace{h_k(x_i)}_{\Delta s_i}$$

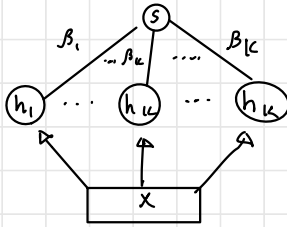
$$s_i = \hat{s}_i + \Delta s_i$$

sequential addition
Epicure

$$= \text{const} + \frac{1}{2} \sum_{i=1}^n \hat{w}_i (\hat{y}_i - h_k(x_i))^2$$

• Root : adaboost

$h_k(x) \in \{+1, -1\}$
weak classifiers



	1 ... j ... p	
1		
⋮	x_i^T	$y_i \in \{+1, -1\}$
n		

→ Theoretical Q:

weak learner = strong learner?

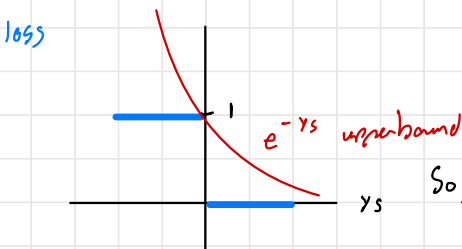
$$s = f(x) = \sum_{k=1}^k \beta_k h_k(x)$$

voting right
p

(v.s XGB, $h_k(x)$ = regression tree = continuous output learned by weighted least squares)

$$y = \text{sign}(s) = \begin{cases} +1 & s \geq 0 \\ -1 & s < 0 \end{cases}$$

• Exponential Loss Function :



if $y_s \gg 0$ then we are very confident in the classification.

So, loss is $\mathbb{1}(y \neq \text{sign}(s)) \leq e^{-y_s}$

continuous convex
encourage big margin y_s

• So our loss function will be:

$$\text{Loss} = \sum_{i=1}^n e^{-y_i s_i}$$

boosting:

$$s = \sum_{k=1}^{K-1} \beta_k h_k(x_i) + \beta_K h_K(x_i)$$

↑ ↑ to be learned
↓ frozen
 $\hat{s}_i + \Delta s_i$

$$\text{Thus } \text{Loss} = \sum_{i=1}^n e^{-y_i (\hat{s}_i + \Delta s_i)} = \sum_{i=1}^n \underbrace{e^{-y_i \hat{s}_i}}_{D_i} e^{-y_i \Delta s_i}$$

$\leftarrow \frac{D_i}{\sum_{i=1}^n D_i}$

$$\sum_{i=1}^n D_i = 1 \quad \text{Distribution Attention}$$

If $-y_i \hat{s}_i$ is big, i receives more attention

- Challenging examples:

$$\text{sign}(y_i) = +1 \quad \text{but} \quad \hat{s}_i < 0$$

$$\text{sign}(y_i) = -1 \quad \text{but} \quad \hat{s}_i > 0$$

$$\text{Loss} = \sum_{i=1}^n D_i e^{-y_i \Delta s_i} = \sum_{i=1}^n D_i e^{-y_i \beta_K h_K(x_i)}$$

$$= \sum_{i=1}^n D_i e^{-\beta_K \underbrace{(y_i h_K(x_i))}_{\substack{+/- \\ +1 \text{ if } y_i = h_K(x) \\ -1 \text{ if } y_i \neq h_K(x)}}}$$

$$= e^{-\beta_K} \sum_{i: y_i = h_K(x_i)} D_i + e^{\beta_K} \sum_{i: y_i \neq h_K(x_i)} D_i$$

$\underbrace{\hspace{10em}}_{1 - \epsilon}$
 $\underbrace{\hspace{10em}}_{\epsilon \text{ error rate}}$

• Recall: $a + b \geq 2\sqrt{ab}$ because $(\sqrt{a} - \sqrt{b})^2 \geq 0$

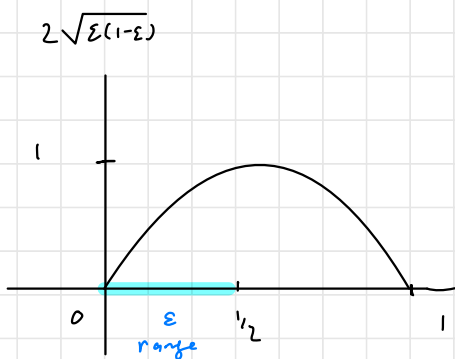
$$\text{Loss} \geq 2\sqrt{e^{-\beta t} \sum_{i: y_i = h_t(x_i)} D_i \cdot e^{\beta t} \sum_{i: y_i \neq h_t(x_i)} D_i} = 2\sqrt{\varepsilon(1-\varepsilon)}$$

$$\text{min: } e^{-\beta t} (1-\varepsilon) = e^{\beta t} \varepsilon$$

$$\frac{1-\varepsilon}{\varepsilon} = e^{2\beta t}$$

$$\text{So } \hat{\beta t} = \frac{1}{2} \log \frac{1-\varepsilon}{\varepsilon}$$

adaptive



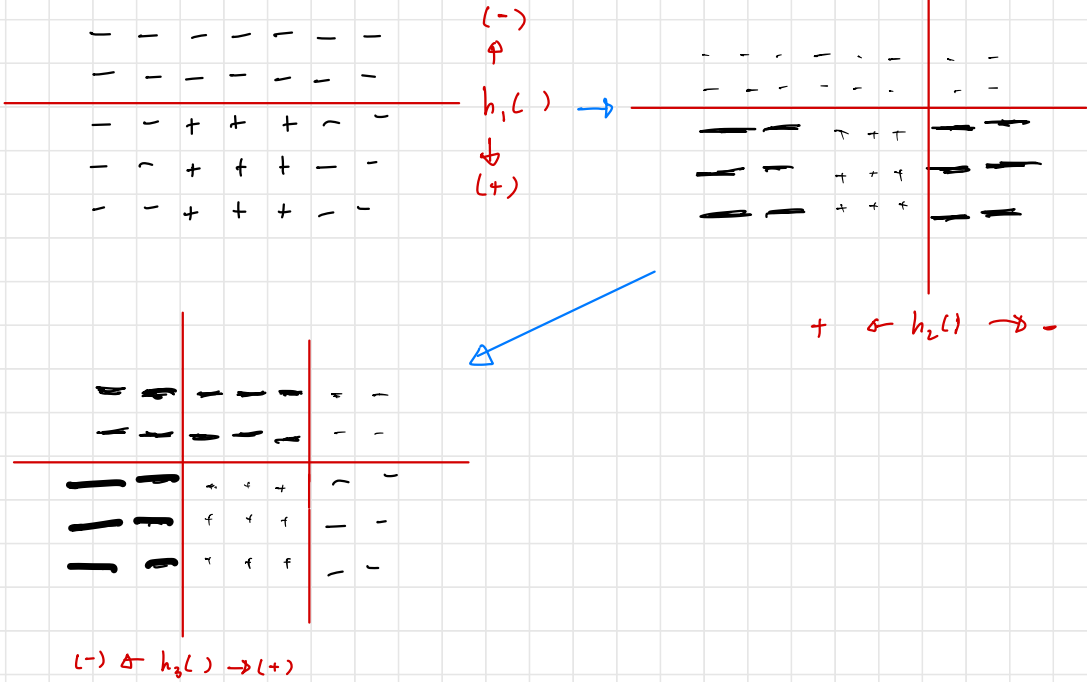
We choose $h_t(\cdot)$ to minimize ε

$$\varepsilon = \sum_{i: y_i \neq h(x_i)} D_i$$

Loss guiding
recursive partitioning



Weak classifier : one layer tree



keep adding trees, tends not to overfit
 increasing margins
 smoothing boundaries

XGB vs Adaboost

Regression trees $\in \mathbb{R}$
 general loss function
 2nd order Taylor
 fit to $\hat{\epsilon}_i$ with \hat{w}_i

classification trees $\in \{+, -\}$
 exponential loss function
 closed form
 $D_i \propto e^{-y_i \hat{w}_i}$

learn from error
 error \rightarrow new tree vs error back-prop

