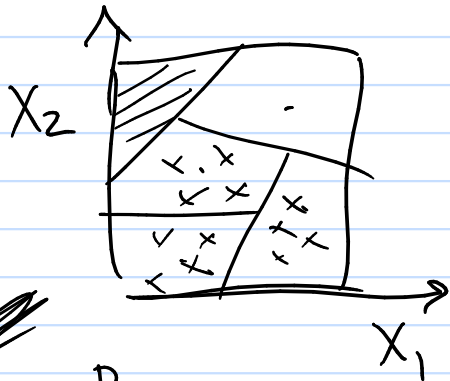Tree model: | Partitioning the input space.
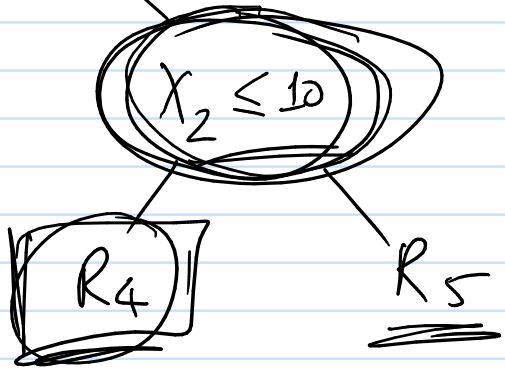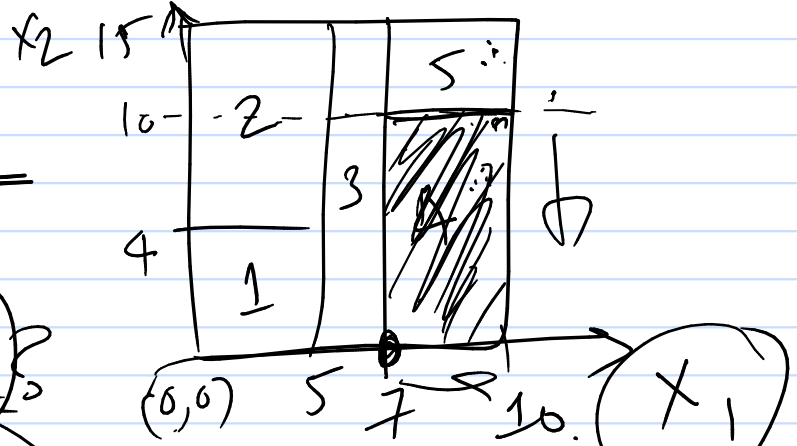
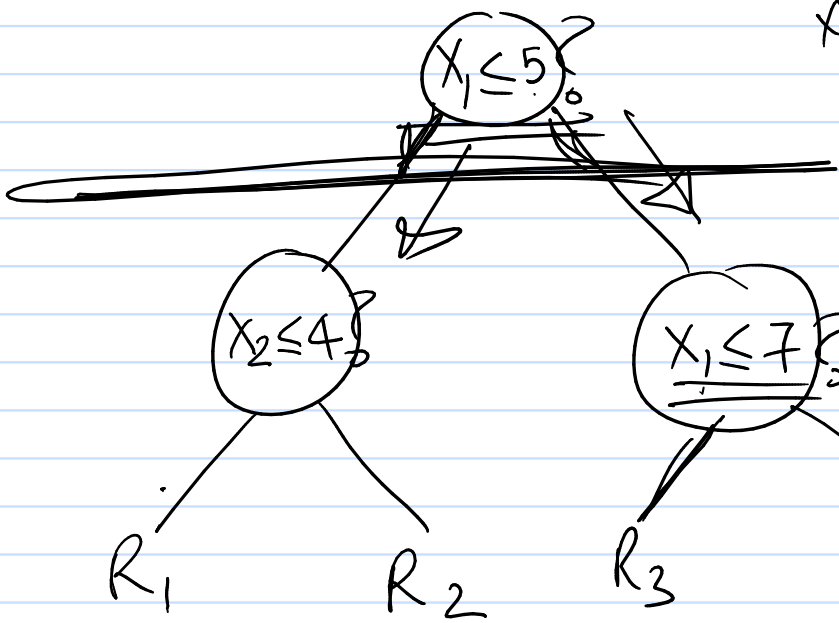$$f(x) = \sum_{j=1}^{J} \widehat{c_j} \, \underline{\mathbb{1}} \left[\underline{x} \in \boxed{R_j}\right]$$

$$x \in \mathbb{R}^P$$

Regular.

$X_1 \leq 5$

$X_2 \leq 4$

$X_1 \leq 7$

$R_1$

$R_2$

$R_3$

$X_2 \leq 10$

$R_4$

$R_5$

$X_2$ 15

10 — 2 —

5

3

4

1

$X_4$

(0,0)   5   7   10.

$X_1$

$$\{R_j, \xi_j\}_{j=1}^{J} = \theta$$

$$\begin{cases} \min_{\theta} \dfrac{1}{N} \sum_{i=1}^{N} \left( y_i - f(x_i) \right)^2 \\ \\ \underbrace{\sum_{j=1}^{J} c_j \, \mathbb{1}\left[ x_i \in R_j \right]} \end{cases}$$

## CART

① Regression tree

$\overline{X_{k}}, \ S$

$X_k \leq S$

$$\dfrac{X_1 \cdots X_p}{\underset{X_k}{\zeta}}$$

$x_2 \leq 10$



$x_1$

$1, 5 = 10$

$$\min_{C_4} \sum_{x_i \in R_4(S)} \left(y_i - C_4\right)^2 + \min_{C_5} \sum_{x_i \in R_5(S)} \left(y_i - C_5\right)^2 =$$
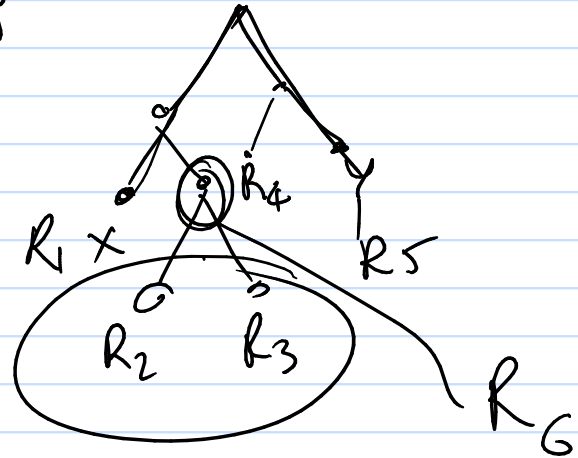
$$\frac{d}{db}(a-b)^2 = -2(a-b)$$

$$-\sum_{x_i \in R_4} 2(y_i - C_4) = 0$$

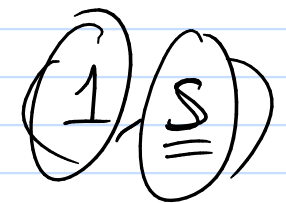$$\frac{1}{\#}\sum y_i = C_4$$

Post processing g a tree : pruning

$$\sum_{\substack{i: \\ x_i \in R_j}} \left( y_i - \boxed{c_j} \right)^2$$

$\underbrace{\hspace{4cm}}$

Node impurity function.

$R_j:$   $\hat{P}_1 = \dfrac{\#\{y_i = 1\}}{\# R_j}$

$\vdots$

$\hat{P}_3$

$K(j) = \underset{\lambda}{\overset{arg}{max}}\left(\hat{P}_1, \hat{P}_2, \hat{P}_3\right)$

Score

$\left(\!\!\begin{array}{c}1\end{array}\!\!\right)\left(\!\!\begin{array}{c}S\end{array}\!\!\right)$ $\longleftarrow$ $\begin{array}{c} R_{j_1} \\ R_{j_2} \end{array}$

$$1 - \hat{P}_{j_1 K(j_1)}$$

$$1 - \hat{P}_{j_2 K(j_2)}$$

$R$

$X_2 < \beta_0$

$R_4$

Gini Index

Gross entropy

$$0.5$$

$$0 \leq \hat{P}_1 \leq 1$$

$$0 \leq \hat{P}_2 \leq 1$$

$$\hat{P}_1 + \hat{P}_2 = 1$$

$$\hat{P}_1$$

$$0 \qquad 1$$

## Issues

① Categorical features

d  f  ©

d  f  c

d  f  c

②    $X_1 \leq S$

$$\sum_{k=1}^{p} \alpha_k X_k \leq S$$

$R_5$

$R_4$

$R_5$

$R_4$

③ Stability / Robustness

④ Smoothness

# Missing data. ∅ // Corruptions / noise

$$X = \begin{bmatrix} & & \text{na} \\ \text{na} & & \\ & \text{na} & \end{bmatrix} \quad Y = \begin{bmatrix} \\ \text{na} \\ \end{bmatrix}$$

$$N \times p \qquad N \times 1$$

## Selection bias

① **Missing at random**

$$P(R \mid X, Y) = P\left(R \mid X^{= \text{observed}}, Y^{\text{observed}}\right) \quad R = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$

② **Missing completely at random.**
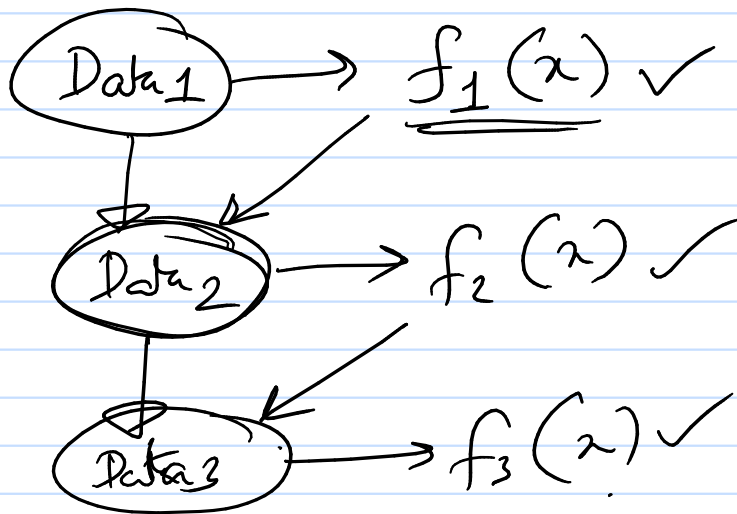
$$P(R \mid X, Y) = P(R)$$

# Forward Stagewise Additive Models.

① Adaboost. M1 $\longrightarrow$ $f(x) = \text{Sign}\left( \sum_{m=1}^{M} \alpha_m f_m(x) \right)$

Binary Classification $+1, -1$

weak
classifiers $\Theta_m$



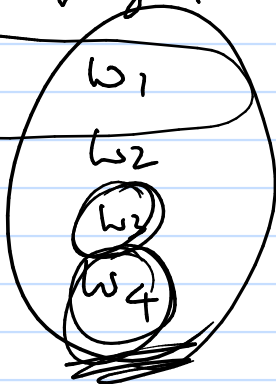Data 1 $\longrightarrow$ $f_1(x)$ ✓

Data 2 $\longrightarrow$ $f_2(x)$ ✓

Data 3 $\longrightarrow$ $f_3(x)$ ✓

$p = 3$

$$X = \begin{bmatrix} 1 & 2 & .5 \\ -1 & .3 & 2 \\ 0.6 & 1.3 & 5 \\ .4 & 1 & 2 \end{bmatrix}$$

weights:

$w_1$

$w_2$

$w_3$

$w_4$

$$Y = \begin{bmatrix} 10 \\ 5 \\ 4 \\ 1 \end{bmatrix}$$

$$\sum_{i=1}^{N} \left( f_2(x_i) - y_i \right)^2 w_i$$

# Adaboost

$$w_i \leftarrow w_i \exp(\alpha_1) \quad \text{if } i \text{ was misClassified}$$

$$\alpha_1 = \text{depends on} \quad \frac{\text{overall}}{\text{misclassification error}}$$

1. FSAM

2. Adaboost $\in$ FSAM

   Why does it work?

3. GBDT $\in$ FSAM

FSAM.

1. $f_1(x) = 0$ ←

for $m = 1, \ldots M:$

(a) $\underline{\theta_m} = \underset{\theta}{\text{argmin}} \sum_{i=1}^{N} \underbrace{(L)} \left( \underline{y_i}, \underbrace{f_{m-1}(x_i)} + b(x_i; \theta) \right)$

(b) $f_m(x) = f_{m-1}(x) + b(x; \theta_m)$

Adaboost $\in$ FSAM $\approx$ approximate.

Binary $\{-1\}$

Loss function $L(y, f(x)) = e^{-\overset{b}{y} f(x)}$ in FSAM
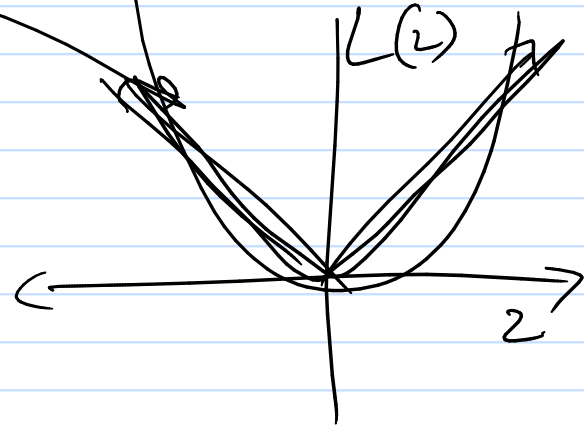
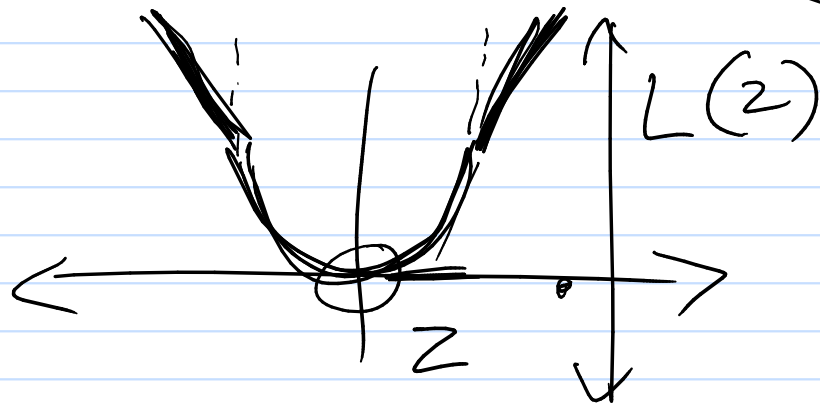$y_i, \quad f_{m-1}(x_i), \quad b(x_i; \theta)$

$\exp(-y_i(f_{m-1}(x_i) + b(x_i; \theta)))$ $\approx w_i$

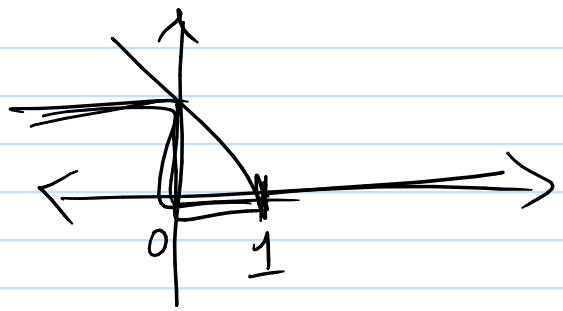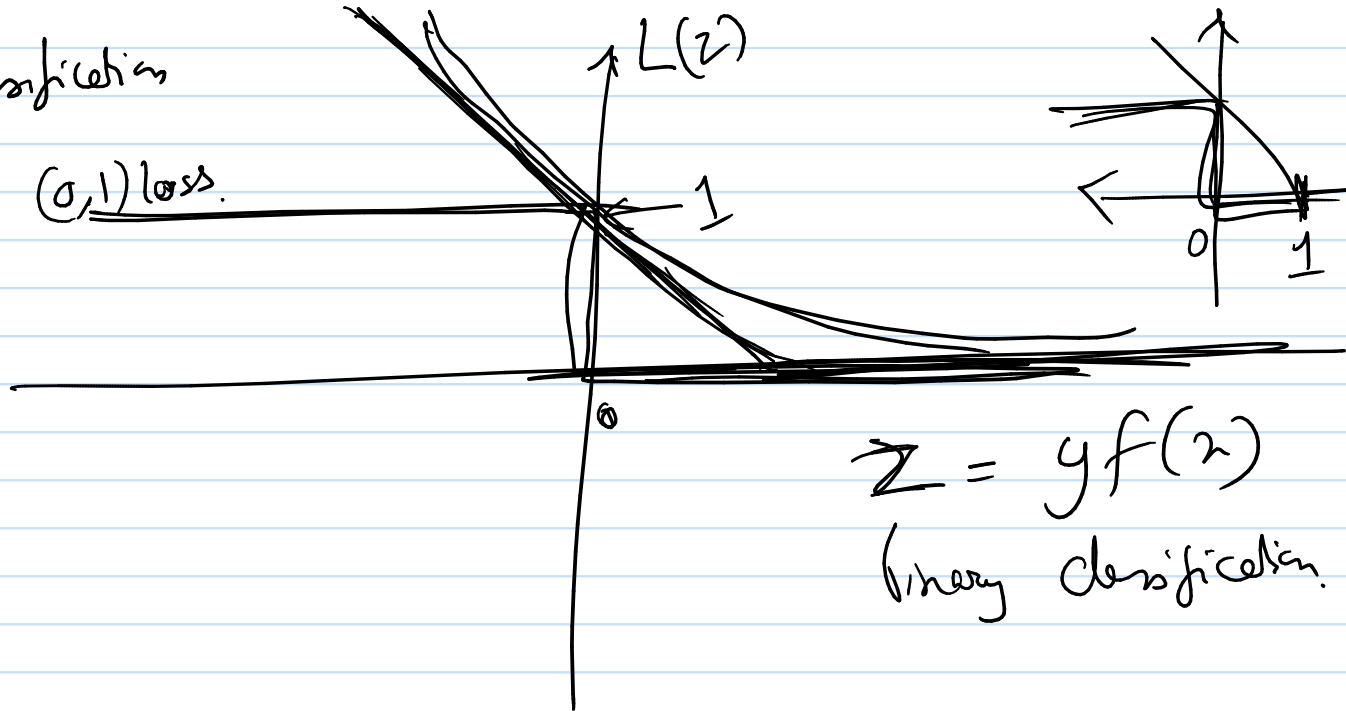$= \boxed{\exp(-y_i f_{m-1}(x_i))} \cdot \exp(-y_i b(x_i; \theta))$

$$L(y, f(x)) = \begin{cases} (y - f(x))^2 & = 2 \quad \text{if } f(x) \text{ is small} \\ |y - f(x)| & \text{if } f(x) \text{ is large.} \end{cases}$$

Huber loss
functin



$L(2)$

$z$

$L(2)$

$z$

Classification

(0,1) loss.



$L(z)$

1

0

$z = yf(x)$

binary classification.

# GBDT :

$f_{m-1}(x)$

$\{x_i, y_i\}_N$
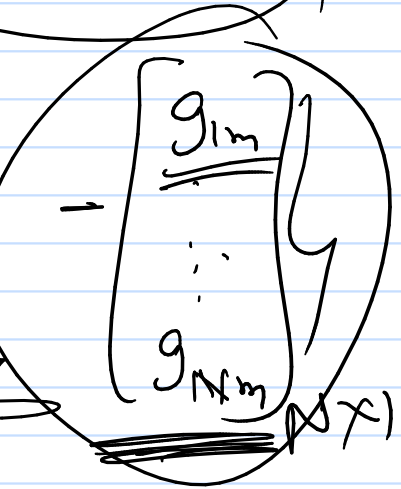
$$g_{im} = \frac{\partial L(y_i, z)}{\partial z} \Bigg|_{z = f_{m-1}(x_i)}$$



$$X \qquad Y_{N \times 1} \qquad \begin{bmatrix} g_{1m} \\ \vdots \\ g_{Nm} \end{bmatrix}(x)$$

$$f_m(x) = f_{m-1}(x) + T(x; \Theta_m)$$

$$\frac{\partial}{\partial z}(y_i - z)^2 = -2(y_i - z)\Big|_{z = f_{m-1}(x_i)} = g_{io}$$

$$= 10$$

$$\min_{\theta} \sum_{i=1}^{N} L\left(y_i, f_{m-1}(x_i) + \underbrace{b(x_i; \theta)}_{z_i}\right)$$

$$\min_{z_1, \ldots z_n} L\left(y_i, \underline{\underline{f_{m-1}(x_i)}} + \boxed{z_i}\right)$$

$$-\frac{\partial L(y_i, z_i)}{\partial z_i}\bigg|_{z_i = f_{m-1}(x_i)} = -g_{im}$$

$$\boxed{f_m(x_i)} = \underline{\underline{f_{m-1}(x_i)}} - \frac{\partial L(y_i, z_i)}{\partial z_i}\bigg|_{z_i = f_{m-1}(x_i)}$$