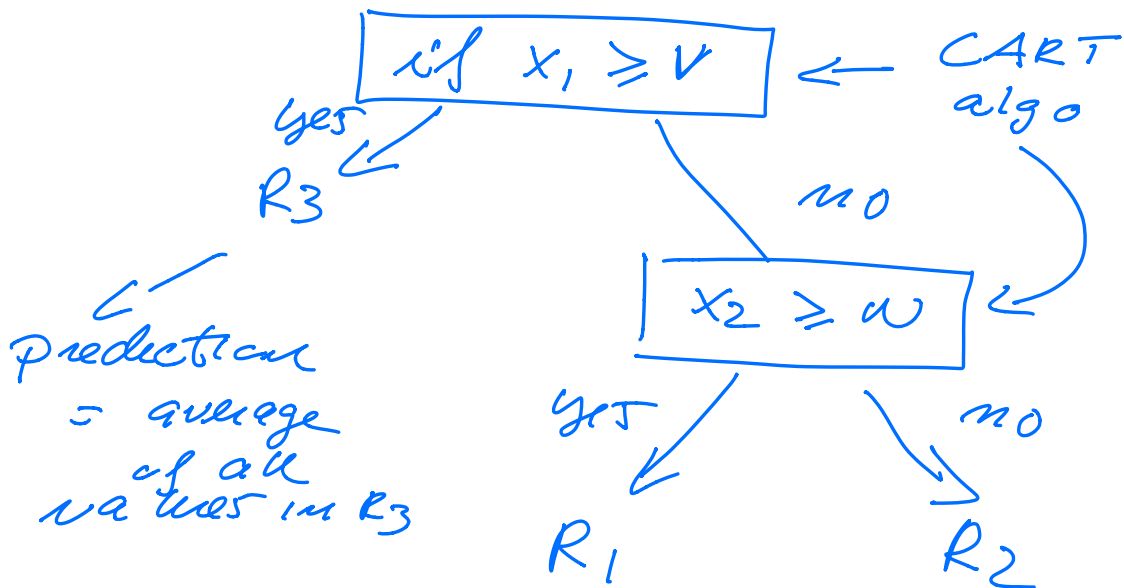
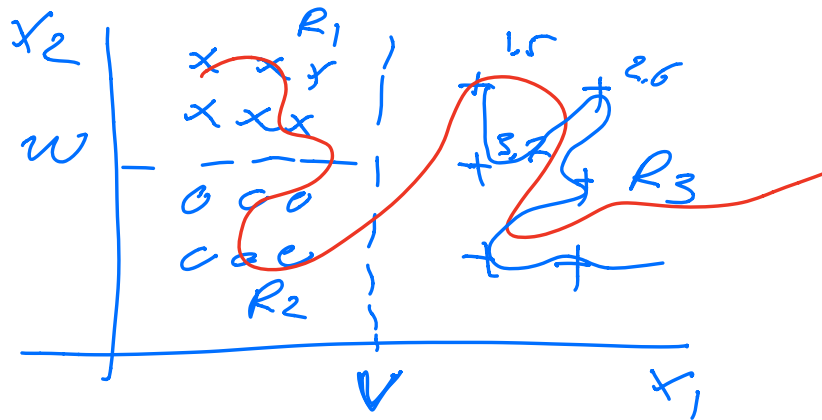
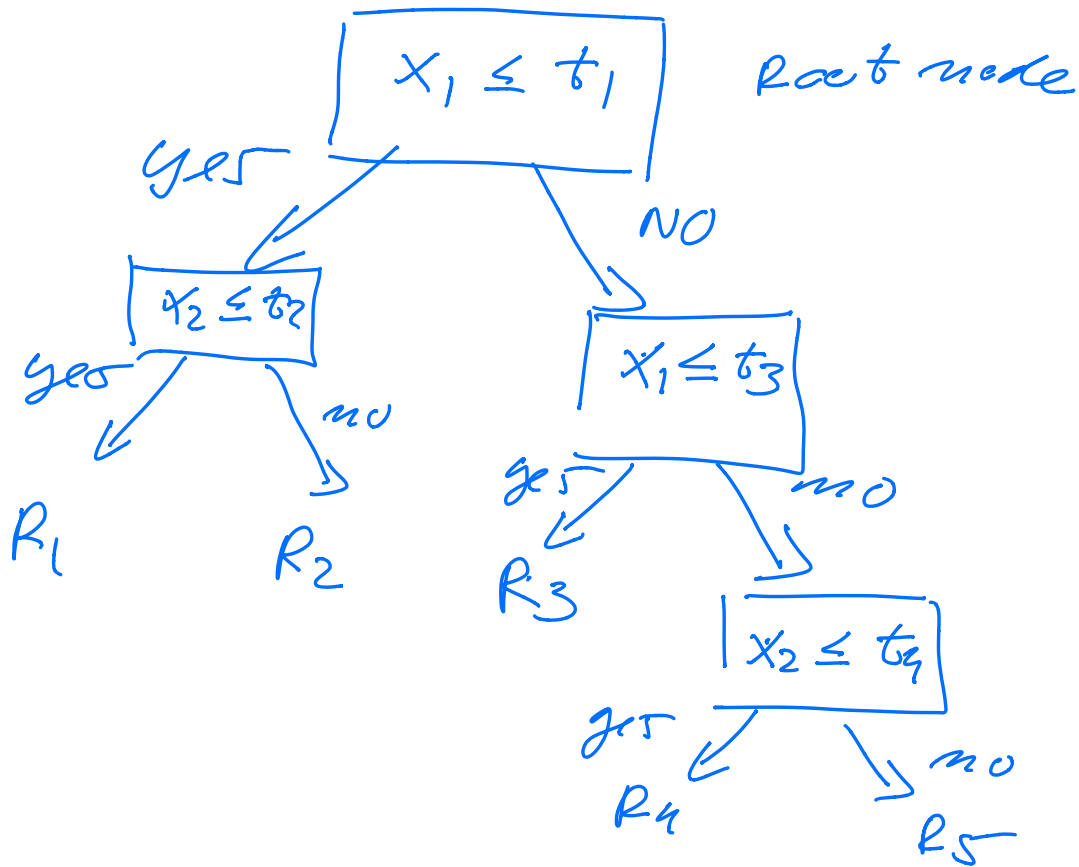
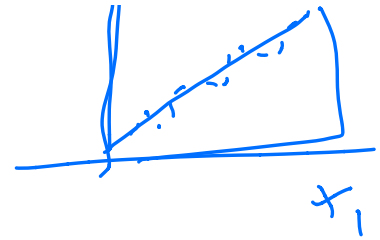
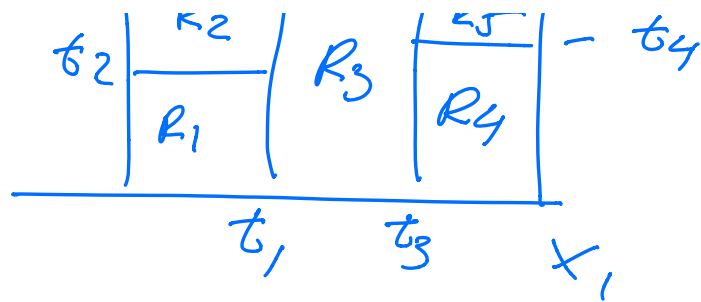


# Lecture January 28



prediction is the mean value of all data points in  $R_1$





Algorithm for Regression case (CART)

$(y_i, x_i) \quad i = 0, 1, \dots, n-1$

$x_i = (x_{i0}, x_{i1}, \dots, x_{ip-1})$

- partition data into  
M-Regions

$R_1, R_2, \dots, R_M$

- Model response as a  
constant  $c_m$  in each  
region

$$f(x) = \sum_{m=1}^M c_m I(x \in R_m)$$

$$\hat{c}_m = \text{ave}(y_i | x_i \in R_m)$$

- Minimize

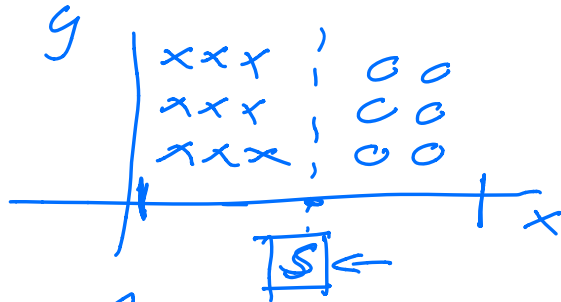
$$\sum_i (y_i - f(x_i))^2$$

$$- R_1(s) = \{x | x_j \leq s\}$$

$$- R_2(s) = \{x | x_j > s\}$$

- Minimize

$$\begin{aligned} & \min_{s} \left[ \min_{c_1} \sum_{x_i \in R_1} (y_i - c_1)^2 \right. \\ & \left. + \min_{c_2} \sum_{x_i \in R_2} (y_i - c_2)^2 \right] \\ & = I(s) \end{aligned}$$



$$\hat{C}_1 = \text{ave}(y_i | x_i \in R_1(s))$$

$$\hat{C}_2 = \text{ave}(y_i | x_i \in R_2(s))$$

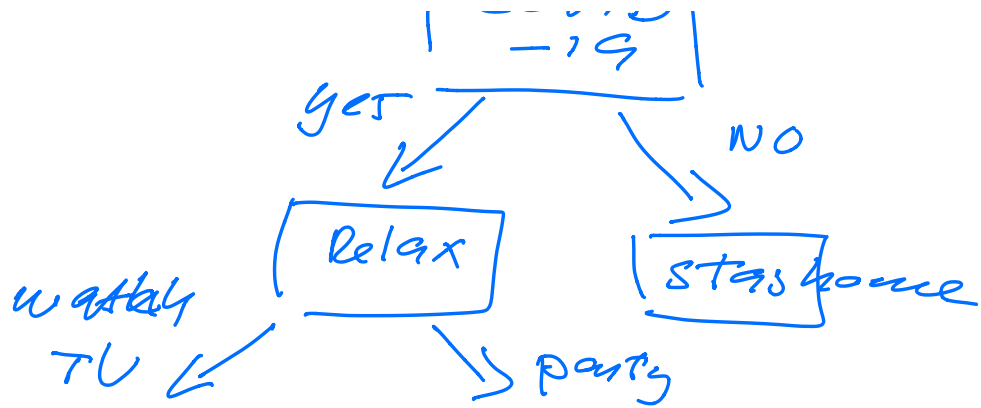
$$J(s) = \frac{N_{\text{left}} \text{MSE}_{\text{left}}}{N} + \frac{N_{\text{right}} \text{MSE}_{\text{right}}}{N}$$

Classification algo (CART)

$$J(s) = \frac{N_{\text{left}} G_{\text{left}}}{N} + \frac{N_{\text{right}} G_{\text{right}}}{N}$$

$G_{\text{left}}/G_{\text{right}}$ ? Gini





Define a node -  $m$  -  
 (Represents a region  $R_m$   
 with  $N_m$  observations)

Proportion of class/category  
 -  $k$  - observations in  
 node  $m$

$$P_{mk} = \frac{1}{N_m} \sum_{x_i \in R_m} I(y_i = k)$$

Classify observations  
 in node -  $m$  -

$$k(m) = \arg \max_k P_{mk}$$

Misclassification

$$(i) \quad \frac{1}{N_m} \sum_{i \in K_m} I(y_i \neq k(m))$$

$$1 - P_{mk}$$

(ii) Gini index (Default in scikit-learn)  
 $K = \text{classes}$

$$\sum_{k=1} P_{mk} (1 - P_{mk})$$

(iii) Entropy

$$- \sum_{k=1}^K P_{mk} \log P_{mk}$$