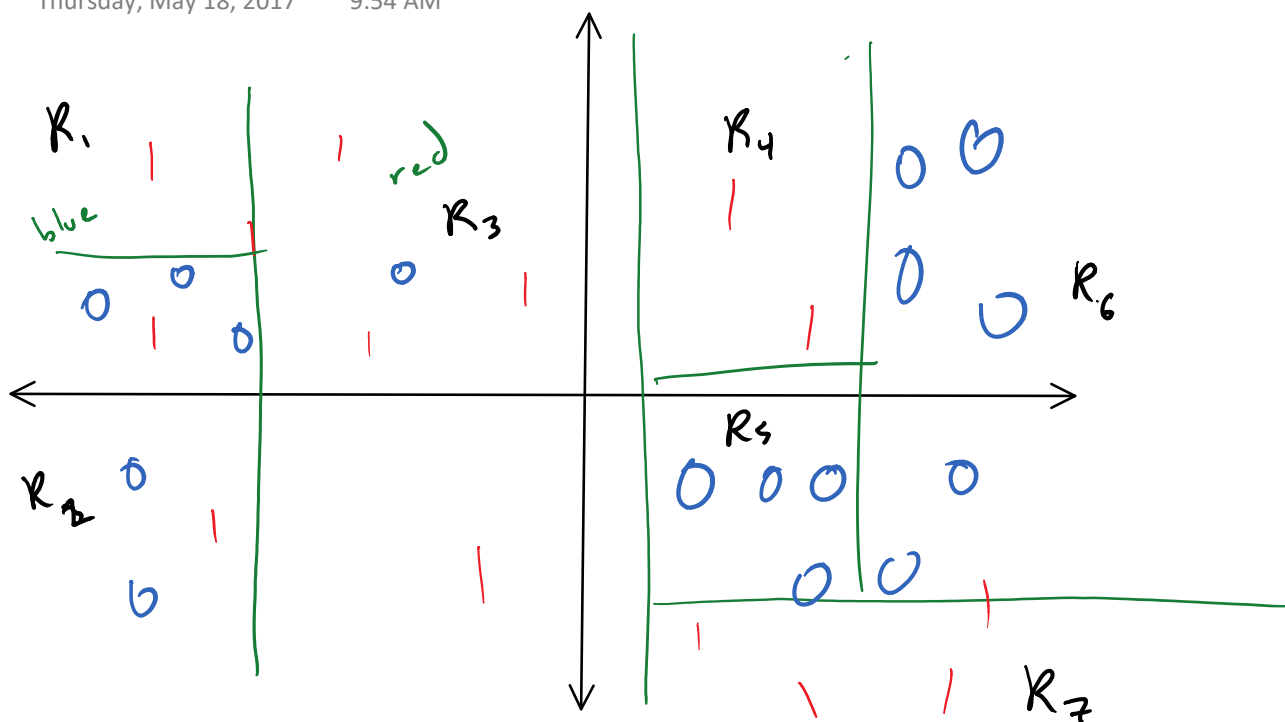


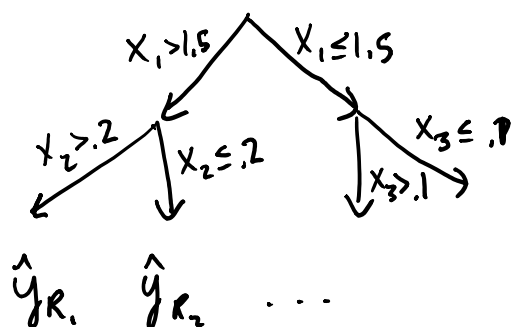
# Decision Trees

Thursday, May 18, 2017

9:54 AM



ex covariates  $X_1, X_2, X_3$



$$\hat{f}(x) = \sum_{m=1}^r \hat{y}_{R_m} \mathbb{1}\{x \in R_m\}$$

(1) Split into halfplanes  $R_1(j, s) = \{x : x_j \leq s\}$   $R_2(j, s) = \{x : x_j > s\}$

$$(2) \min_{j, s} \left[ \min_{\hat{y}_1} \sum_{x_i \in R_1(j, s)} l(y_i, \hat{y}_1) + \min_{\hat{y}_2} \sum_{x_i \in R_2(j, s)} l(y_i, \hat{y}_2) \right] (*)$$

(3)  $y_i \leftarrow y_i - \hat{y}_{1/2}$  ← depend on region

(4) repeat (1)-(3) w/  $R_1, R_2$

Regression:  $l(y_i, \hat{y}) = (y_i - \hat{y})^2$       $\hat{y}_i = \frac{1}{|R_{(i,s)}|} \sum_{x_i \in R_{(i,s)}} y_i$

Classification: classes  $k \in \{1, \dots, K\}$

▷  $l(y_i, \hat{y}) = \mathbb{1}\{y_i \neq \hat{y}\}$  (not use (3))

w/ rectangle pred.  $R_m$       $\hat{p}_{mk} = \frac{1}{|R_m|} \sum_{x_i \in R_m} \mathbb{1}\{y_i = k\}$

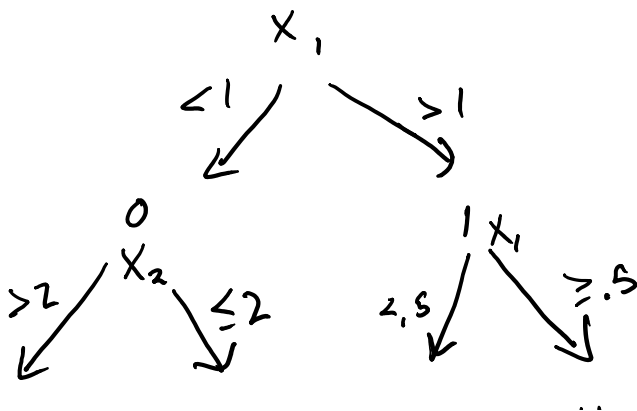
$\hat{y}_i = \arg \max_k \{\hat{p}_{mk}, x_i \in R_m\} \quad (x_i \in R_m)$

▷ Gini Index      $\sum_{k \neq k'} \hat{p}_{mk} \hat{p}_{mk'} = \sum_k \hat{p}_{mk} (1 - \hat{p}_{mk})$

▷ Cross-Entropy      $-\sum_k \hat{p}_{mk} \log \hat{p}_{mk}$

— Gini & CE are differentiable, prefer "balanced" trees, outputs "soft classifier"

Code: Binary string - digits represent splits



Length of codeword is depth of tree

MDL principle use  $L$  as a tuning parameter

00

01

10

11

AIC:  $-\log\text{lik} + L$