

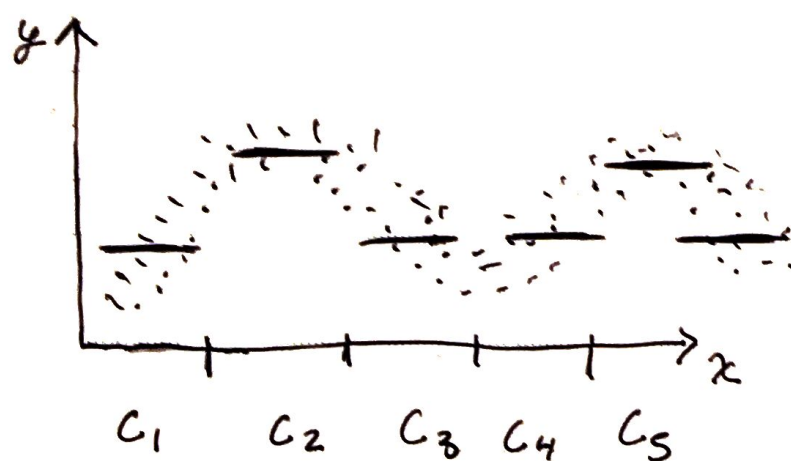
## Math 342W Lecture 18

### Classification and Regression Trees

Classification Trees:  $y = \{C_1, C_2, \dots, C_n\}$

Regression Trees:  $y \in \mathbb{R}$

Lets consider  $p\text{-row} = 1$  and build up Regression Trees...



$$f: \mathbb{R}^p \rightarrow \mathbb{R}$$

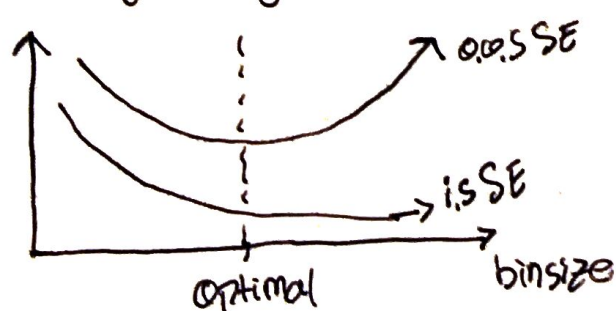
We start by splitting the domain and categorize it, choosing  $\bar{y}$  within each category.

$$\mathcal{H} = \{ \omega_1 \mathbb{1}_{x \leq 1} + \omega_2 \mathbb{1}_{x \in (1,2]} + \dots + \omega_p \mathbb{1}_{x > p} : \vec{\omega} \in \mathbb{R}^p \}$$

What choice of bin size would be best?  $\rightarrow$  Hyperparameter

Large binsize  $\rightarrow$  Underfit and Small binsize  $\rightarrow$  Overfit

We can use the model selection procedure to optimize this by using ~~cos~~ metrics.



With equally sized bins the total bins are defined by  $B^p$ , which is uncontrollably large.

$\hookrightarrow$  Change bins based on  $\frac{d}{dx}$

We can restrict our bin splitting to be orthogonal to axis. In this way we can impose limits on the algorithm

## Regression Trees Algorithm

- ① Begin with  $D = \langle X, \vec{y} \rangle$
- ② Consider all possible orthogonal to axis splits  $\langle X_L, \vec{y}_L \rangle, \langle X_R, \vec{y}_R \rangle$  with rules

$$\left. \begin{array}{l} x_1 \leq x_{11}, x_1 \leq x_{12}, \dots, x_1 \leq x_{1(n-1)} \\ x_2 \leq x_{21}, x_2 \leq x_{22}, \dots, x_2 \leq x_{2(n-1)} \\ \vdots \\ x_p \leq x_{p1}, x_p \leq x_{p2}, \dots, x_p \leq x_{p(n-1)} \end{array} \right\} \begin{array}{l} \text{we only consider} \\ \text{to the } (n-1)^{\text{th}} \\ \text{bin or else it} \\ \text{would be 1 bin.} \end{array}$$

For each split there will be two daughter nodes.  
Compute  $SSE_L$  and  $SSE_R$

- ③ Locate split with lowest weighted average  $\frac{n_L SSE_L + n_R SSE_R}{n_L + n_R}$
- ④ Assign  $\hat{y}$ 's to be  $\bar{y}$ 's of two bins
- ⑤ Recursively repeat ①-④ for daughter bins until bin has  $n_0 \leq \text{Observations}$ , where  $n_0$  is a hyperparameter

## Classification Trees Algorithm

- ① " "
- ② " "  $Gini_L$  &  $Gini_R$
- ③ " "  $\frac{n_L Gini_L + n_R Gini_R}{n_L + n_R}$
- ④ " "  $\text{mode}[y\text{'s}]$

$$\left. \begin{array}{l} Gini_L = \sum_{L=1}^L \hat{P}_L (1 - \hat{P}_L) \\ \hat{P}_L = \frac{\# y_L}{n_L} \end{array} \right\}$$