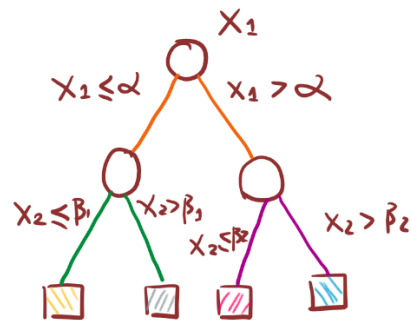
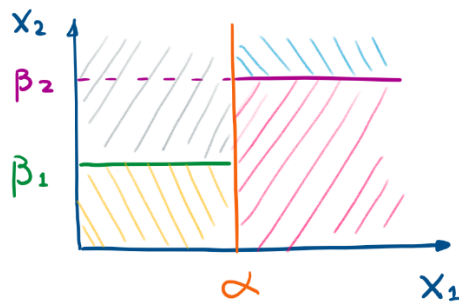


CART Note

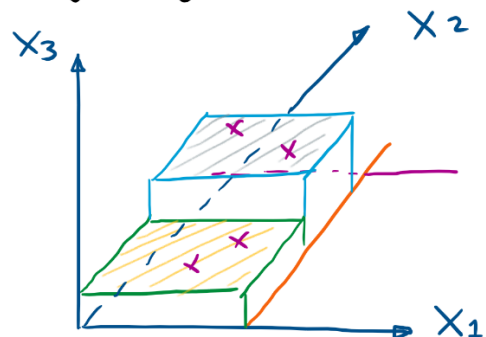
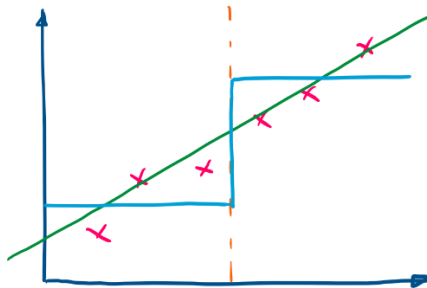
Zengchang Qin

CART - classification and Regression Tree

- 1) We have learned that a decision tree is a model that divides the input space into a few regions.



- 2) If we only consider one-dimensional variable x .
how can we use a tree model for regression?



3) Given training data $D = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(n)}, y^{(n)})\}$
 We hope to learn a CART tree that minimize the following cost function:

$$\text{Loss} = \min_{j, s} [\min_{C_1} L(y^{(i)}, C_1) + \min_{C_2} L(y^{(i)}, C_2)]$$

$$\text{Where } C_m = \text{ave}(y_i | x_i \in R_m)$$



In the CART Model, we need to find the best variable X_j and cut point s in order to minimize the loss function. We can rewrite the loss by:

$$\text{Loss} = \min_{j, s} [\min_{C_1} \sum_{x \in R_1(j, s)} (y^{(i)} - C_1)^2 + \min_{C_2} \sum_{x \in R_2(j, s)} (y^{(i)} - C_2)^2]$$

$$\text{where } R_1(j, s) = \{x | x_j \leq s\}$$

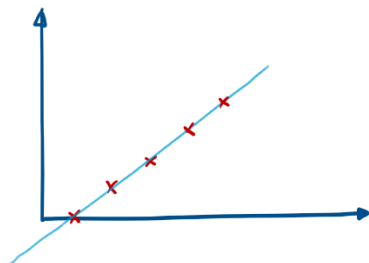
$$R_2(j, s) = \{x | x_j > s\}$$

$$C_m = \frac{1}{N_m} \sum_{x \in R_m(j, s)} y^{(i)}, m = 1, 2$$

4) An example:

x	1	2	3	4	5	6	7	8	9	10
y	0	1	2	3	4	5	6	7	8	9
		1.5			4.5					

($j=1$)



Given $s = 1.5$

$$R_1 = \{0\} \quad C_1 = 0$$

$$R_2 = \{1, 2, 3, \dots, 9\} \quad C_2 = 5$$

$$\text{Loss}(1.5) = 0^2 + \frac{1}{9} \sum_{i=1}^9 (i-5)^2$$

⋮

when $s = 4.5$

$$R_1 = \{0, 1, 2, 3\} \quad C_1 = 1.5$$

$$R_2 = \{4, 5, 6, \dots, 9\} \quad C_2 = 6.5$$

$$\text{Loss}(4.5) = \sum_{i=0}^3 (i-1.5)^2 + \sum_{j=4}^9 (j-6.5)^2$$

We need to calculate $\text{Loss}(s)$ for $s=1.5, 2.5 \dots 9.5$ to find the best cut point s . If we have more than one variable, we need to search ($j=1, 2 \dots n$) with best cut point s_j .