

$$\mathcal{H} = \left\{ w_1 b_1(\vec{x}) + w_2 b_2(\vec{x}) + \dots + w_B b_B(\vec{x}) : \vec{w} \in \mathbb{R}^B \right\} \quad y = \mathbb{R}$$

and b_1, b_2, \dots, b_B are known functions that attempt to span the function space of $f: \mathbb{R}^p \rightarrow \mathbb{R}$.

Example set of functions: set of all first-order interactions

$$1, x_1, x_2, \dots, x_p, x_1^2, x_2^2, \dots, x_p^2, x_1 x_2, x_1 x_3, \dots, x_{p-1} x_p \quad B = 2p + \binom{p}{2}$$

Set of all second order interactions
 $x_1 x_2 x_3, \dots, x_{p-2} x_{p-1} x_p$
 $\sum x_i^{d_i} x_j^{d_j} x_k^{d_k} : i, j, k \text{ unique}, d_1 + d_2 + d_3 \leq 3$
 $d_n \in \{0, 1, 2, 3\}$

B is exponentially large

A: likely... w 's will be sparse i.e. most w 's = 0.

② Let $g(\vec{x}) = g_0(\vec{x}) = \bar{y}$, $X = \begin{bmatrix} \vec{1} \end{bmatrix}$ OLS regression

① Try all B individually. $y \sim b_1(\vec{x}), y \sim b_2(\vec{x}), \dots, y \sim b_B(\vec{x})$ and compute SSE

$$X = \begin{bmatrix} \vec{1} & b_1(\vec{x}) \end{bmatrix}$$

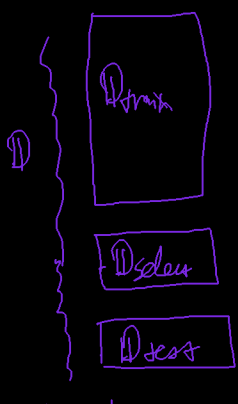
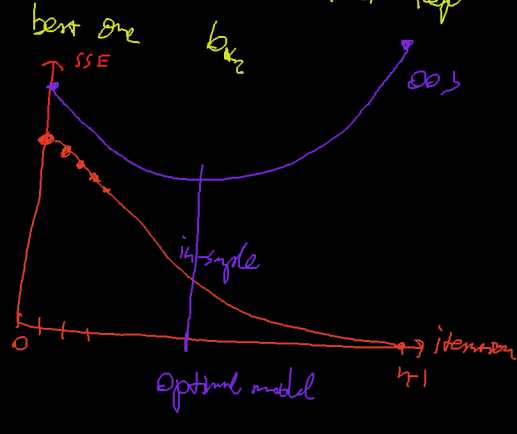
Reduction for each and keep best one b_{k_1}

② Try all $B-1$ remaining individually $y \sim \begin{bmatrix} \vec{1} & b_{k_1}(\vec{x}) & b_{k_2}(\vec{x}) \end{bmatrix}$ and compute SSE

Reduction for each and keep best one b_{k_2}

Repeat

③ Stop if... crossSSE goes up.

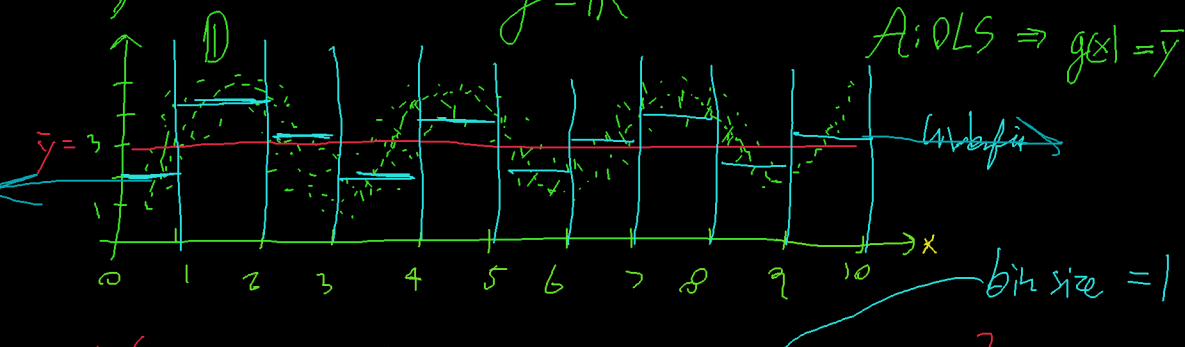


Model Selection Procedure

$$g_1, g_2, \dots, g_m \quad m = 2^B$$

The algorithm could be modified to use K-fold CV (inner & outer).

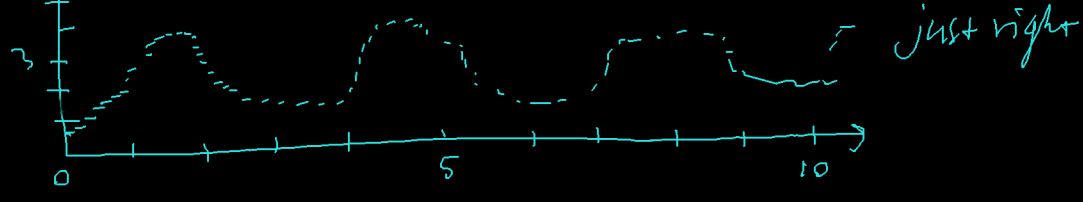
Classification and Regression Tree (CART) Algorithm (1984).



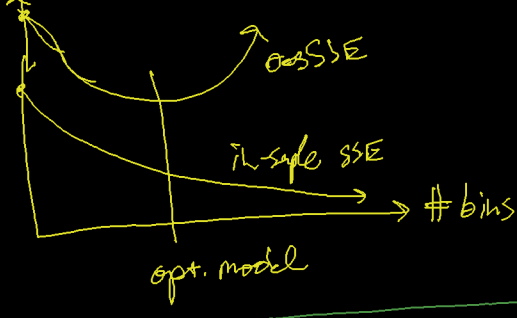
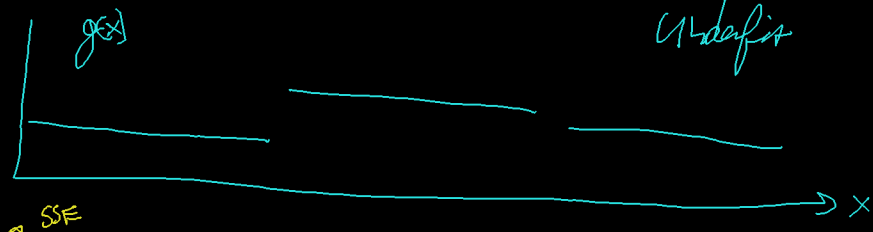
$$\mathcal{H}_B = \left\{ w_1 \mathbb{1}_{x \in [0,1]} + w_2 \mathbb{1}_{x \in (1,2]} + \dots + w_9 \mathbb{1}_{x \in [9,10]} : \vec{w} \in \mathbb{R}^9 \right\}$$

A: OLS $\Rightarrow g(x) =$

$$\mathcal{H}_{B0.1} = \left\{ w_1 \mathbb{1}_{x \in [0,0.1]} + \dots + w_{99} \mathbb{1}_{x \in [9.9,10]} : \vec{w} \in \mathbb{R}^{99} \right\}$$



$$\mathcal{H}_{B3.3} = \left\{ w_1 \mathbb{1}_{x \in [0,3.3]} + w_2 \mathbb{1}_{x \in (3.3,6.7]} + w_3 \mathbb{1}_{x \in (6.7,10]} : \vec{w} \in \mathbb{R}^3 \right\}$$

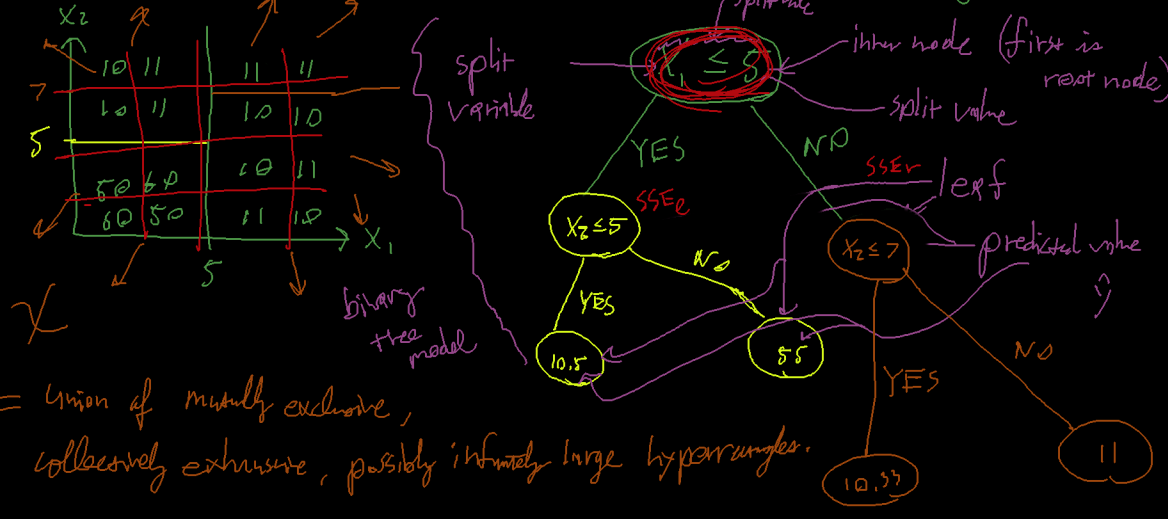


Select bin size using the crossSSE curve

In two dimensions ($p=2$), the bins are squares. $B = \# \text{ bins} / \text{dim} \Rightarrow B_{\text{best}} = B^2$

In p dimensions $B_{\text{best}} = B^p$ which $> n$ very fast.

A solution: "binary trees" with split rules that are orthogonal to an axis.



= Union of mutually exclusive, collectively exhaustive, possibly infinitely large hyperangles.

Regression Tree Algorithm

① Let dataset be all data

ordered / sorted values

② Consider every possible orthogonal-to-axis split $x_j \leq x_{(j)}$ $j=1 \dots p, i \in 1 \dots n-1$. and compute SSE_l, SSE_r the SSE's in the partition left node and right node. Select the rule where

$$SSE_{\text{weighted}} := \frac{n_l SSE_l + n_r SSE_r}{n_l + n_r} \quad n_l := \# \text{ observations in left set} \quad n_r := \# \text{ observations in right set}$$

is smallest, i.e. create an inner node with this split rule and a left leaf with $\hat{y} = \bar{y}_l$ and a right leaf with $\hat{y} = \bar{y}_r$.

③ If $n_l > N_0$ then set dataset = left partition and run step 1 on it. If $n_r > N_0$ then set dataset = right partition and run step 1 on it.

Also: don't split if all the y's are the same

N_0 is a hyperparameter. If N_0 is small e.g. $N_0 = 1 \Rightarrow g$ is overfit. If N_0 large \Rightarrow underfit model. How to pick N_0 ? Use 3-fold selection.