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

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

Example set of function: Set of all first order interactions

$$X_1, X_2, \ldots, X_p, X_1^2, X_2^2, \ldots, X_p^2, X_1 X_2, X_1 X_3, \ldots, X_{p-1} X_p$$

$B = 2p + \binom{p}{2} =$ ,, ,, ,, ,, ,, ,, ,, $X_1^3$,

$X_2^3, \ldots, X_p^3, X_1^2 X_2, \ldots, X_{p-1}^2 X_p, X_1 X_2 X_3, \ldots, X_{p-2} X_{p-1} X_p$

Set of all second order interactions.
B is exponensially large.

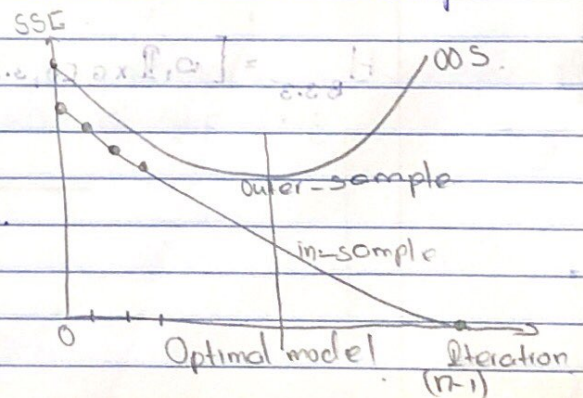A: Likely ... W's will be spark i.e. most w's = 0.

Forward stepwise OLS

⓪ Let $g(\vec{x}) = g_o(\vec{x}) = \bar{y}$, $X = [\vec{1}]$       OLS regression

① Try all B individually. $y \sim b_1(\vec{x})$, $y \sim b_2(\vec{x}), \ldots, y \sim b_B(\vec{x})$
and compute SSE reduction for each $X = [\vec{1} \; b_k(\vec{x})]$.
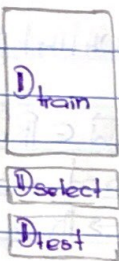and keep best one $b_{*_1}$

② Try all B-1 remaining individually; $y \sim [\vec{1} \; b_{*_1}[\vec{x}] \; b_k[\vec{x}]]$
and update SSE reduction for each and keep
best one $b_{*_2}$.
repeat.

⋮

⑤ Stop if ... OOS SSE goes up.

(II)



$D_{train}$

$D_{select}$

$D_{test}$
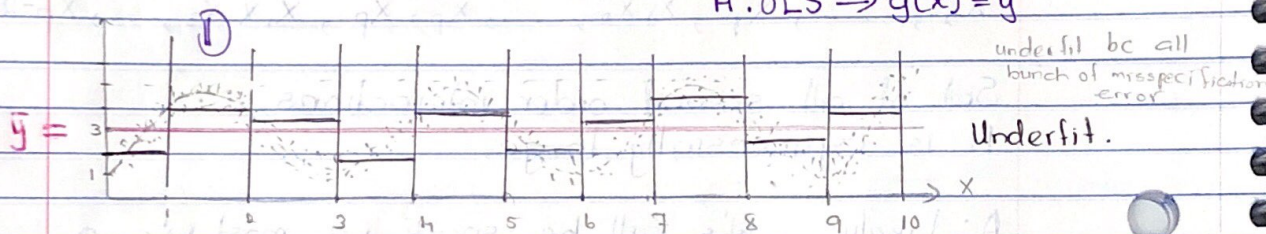
Model selection procedure $g_1, g_2, \ldots, g_m$

$m = 2^B$

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

$y = \mathbb{R}$

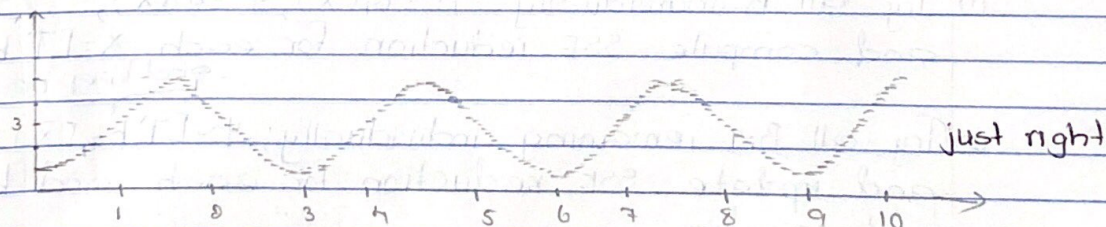## Classification & Regression Tree (CART) Algorithm (1984)

①

$A: OLS \Rightarrow g(x) = \bar{y}$

$\bar{y} = 3$

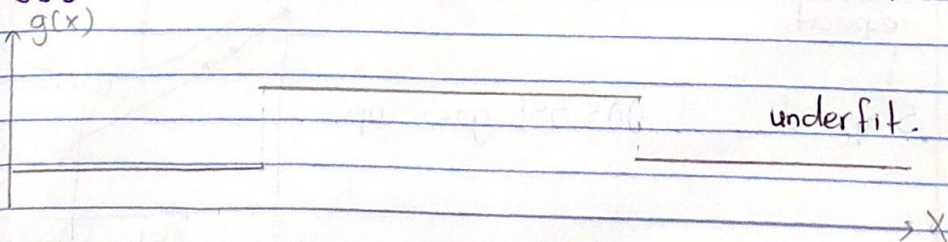underfit bc all bunch of misspecification error

Underfit.

bin size = 1

$$H_{B1} = \left\{ w_1 \mathbb{1}_{x \in [0,1]} + w_2 \mathbb{1}_{x \in [1,2]} + \ldots + w_9 \mathbb{1}_{x \in [9,10]} : \vec{w} \in \mathbb{R}^9 \right\}$$
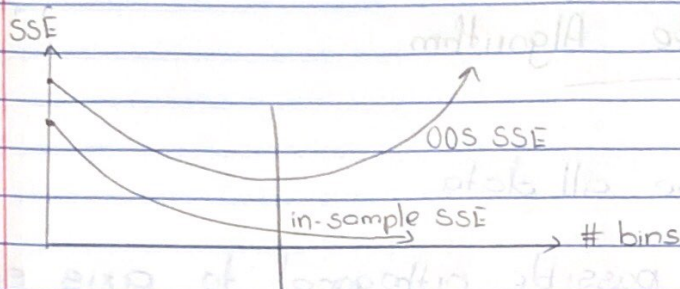
$A: OLS \Rightarrow g(\vec{x}) =$

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

3

just right

$$H_{B3.3} = \left\{ w_1 \mathbb{1}_{x \in [0,3.3]} + w_2 \mathbb{1}_{x \in [3.3,6.6]} + w_3 \mathbb{1}_{x \in [6.6,10]} : \vec{w} \in \mathbb{R}^3 \right\}$$

$g(x)$

underfit.

SSE



OOS SSE

Select bin size using the OOS SSE curve

in-sample SSE

→ # bins

P=1 ↑

In two dimentions $(p=2)$, the "bins" are squares.
$B = \#$ bins/dim $\implies B_{tot} = B^2$

In p dimentions $B_{tot} = B^p$ which $> n$; very fast.

A solution: "Binary trees" with split rules that are othogonal to an axis.



Split variable → $X_1 \leq 5$ ← split rule
inner node (first is root node)
split value

YES          NO
SSE $l$      SSE $r$

$X_2 \leq 5$      $X_2 \leq 7$

binary tree model

| | | | |
|---|---|---|---|
| 10 | 11 | 10 | |
| 10 | 11 | 11 | 10 |
| 50 | 60 | 10 | 11 |
| 60 | 50 | 11 | 10 |

7
5

5

YES  NO      YES  NO

Predicted value → 10.5    55    10.33    11
      ŷ

leaf

= Union of mutually exclusive, collectively exhausser possibly. infinitely large hyperrectangles.

# Regression Tree Algorithm

◎ Let dataset be all data

① Consider every possible orthogonal to axis split
$$X_j < X_{ij} \quad j=1\ldots p, \; IGI, \ldots, n-1 \qquad \uparrow i=q$$
↑
ordered/sorted values

And calculate $SSE_l$, $SSE_r$ the SSE's in the putative left node and right node. Select the rule where.

$$SSE_{weighted} = \frac{n_l \, SSE_l + n_r \, SSE_r}{n_l + n_r}$$

$n_l$: # observations in left set.
$n_r$: # observations in right set
is smallest. i.e create an inner node with that split rule and a left leaf and with $\hat{q} = \bar{y}_l$ and a right leaf with $\hat{q} = \bar{y}_r$.

② 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,
(No is a hypermater. If $N_0$ is small e.g. $N_0 = 1 \Rightarrow$
q is overfit.)
If $N_0$ large $\Rightarrow$ Underfit model.
How to pick $N_0$? Use 3-fold selection.