## Lecture 18:

$$H = \{ W_1 b_1(\vec{X}) + W_2 b_2(\vec{X}) + \cdots + W_B b_B(\vec{X}) ; \vec{W} \in \mathbb{R}^B \}$$

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 functions: - set of all first order
interaction

$$X_1, X_2, \ldots, X_P, X_1^2, X_2^2, \ldots, X_P^2 ; X_1 X_2, X_1 X_3, \ldots, X_{P-1} X_P$$

$$\ldots, X_1^3, X_2^3, \ldots, X_P^3,$$

$$X_1 X_2, \ldots, X_{P-1} X_P,$$

$$X_1 X_2 X_3, \ldots, X_{P-2} X_{P-1} X_P.$$

$$\boxed{B = 2P + \binom{P}{2} \\ = 2^P}$$

Set of all second interactions

$B$ is exponentially large.

A: likely... $W$'s will be <u>sparse</u>  i.e. most $W's = 0$

Forward stepwise OLS
- ⓪ let $g(\vec{X}) = g_0(\vec{X}) = \overline{Y}$  ⟋ OLS represent $\infty$ $n_1$
- 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
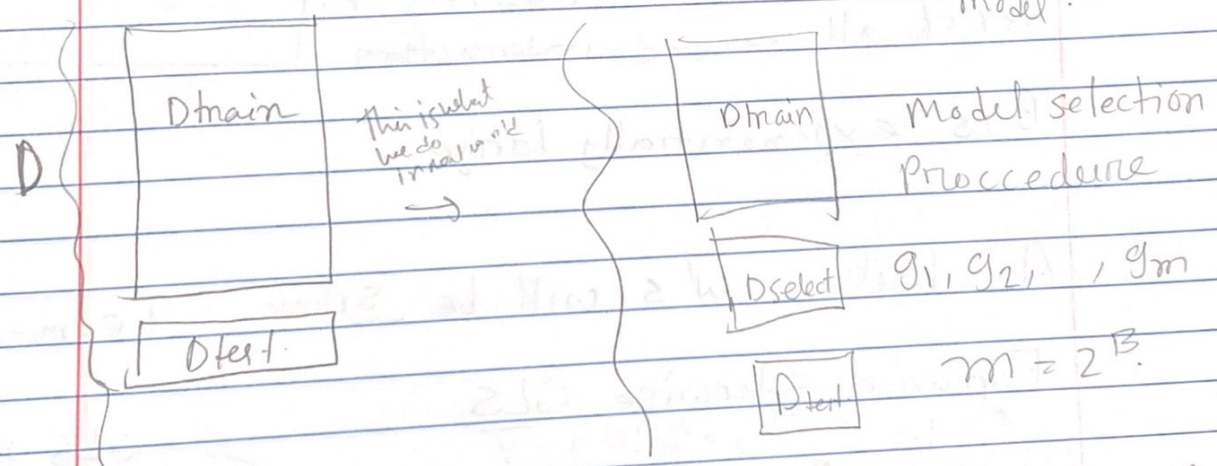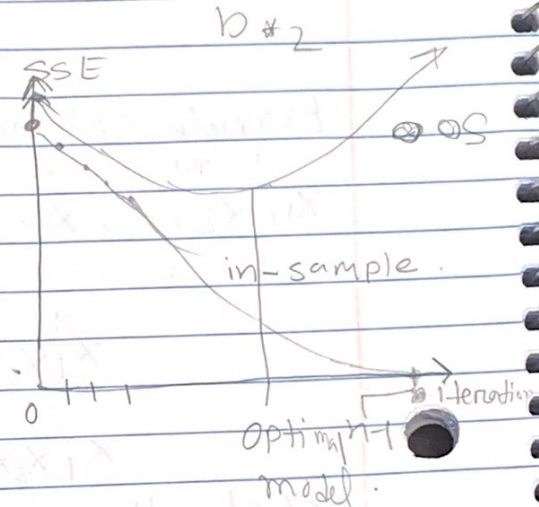  each.
  $$X = \left[ \overline{1} \cdot b_K(\vec{X}) \right]$$

② Try all $B-1$ remaining individually

$y \sim [\vec{1} \ b_{*_1}(\vec{x}) \ b_k(\vec{x})]$ and compute

SSE reduction for each and best one

$b_{*_2}$

:

repeat

③ stop if ... (outer sample) OOSSE goes up.



SSE

in-sample

0 ⟶ iteration

Optim$_{n-1}$ model.

② OS

$D$ ⎰ 
| Dtrain |
this is what we do in real work ⟶

| Dtrain |  Model selection
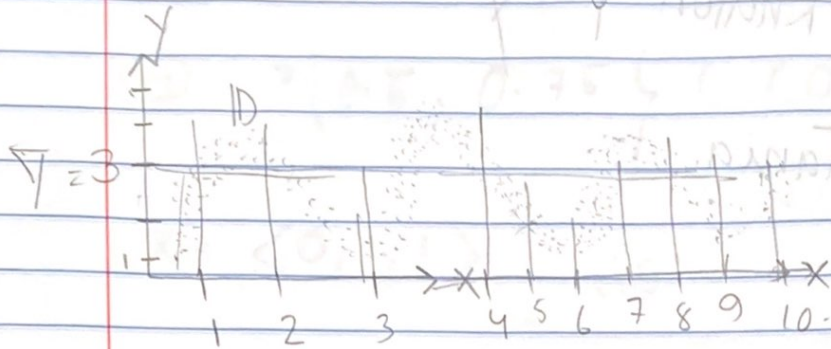| Dselect |  $g_1, g_2, \ldots, g_m$
| Dtest |  $m = 2^B$

| Dtest |

— The algorithm could be modified to use $k$-fold CV (inner & outer)

$$y = \mathbb{R}$$

## Classification and regression Tree CarT Algorithm (1984)
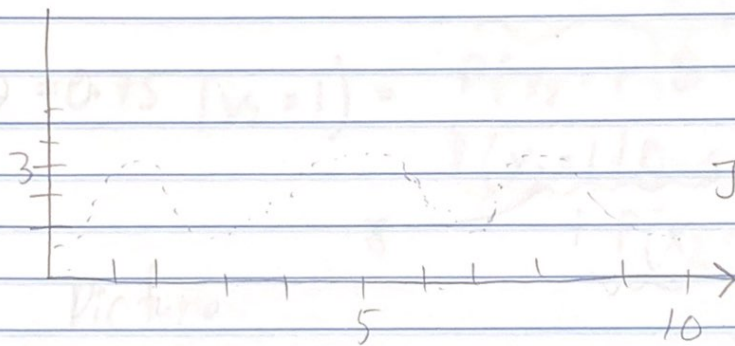
$$A : OLS = g(X) = \overline{Y}$$

underfit

$$H_{B_1} : \{w_1 \mathbb{1}_{x \in [0,1]} + w_2 \mathbb{1}_{x \in [1,2]} + \cdots + w_9 \mathbb{1}_{x \in [9,10]} ; \quad w \in \mathbb{R}^9 \}$$
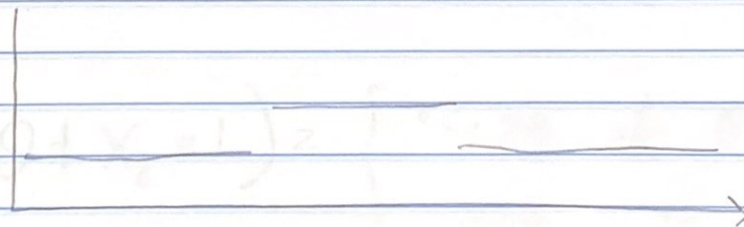
$$A : OLS \Rightarrow g(\overline{X}) = $$

$$H_{B_{0.1}} = \{w_1 \mathbb{1}_{x \in [0,\cdot1]} + \cdots + w_{99} \mathbb{1}_{x \in [9.9,10]} ; \quad 100 \, w \in \mathbb{R}^{99} \}$$

3

Just right

5          10

$$H_{B_{3.3}} = \{w_1 \mathbb{1}_{x \in [0,3\cdot3]} + w_2 \mathbb{1}_{x \in [3.3, 6.6]} $$

underfit.

SSE

OOSSE

in sample SSE

# bins.

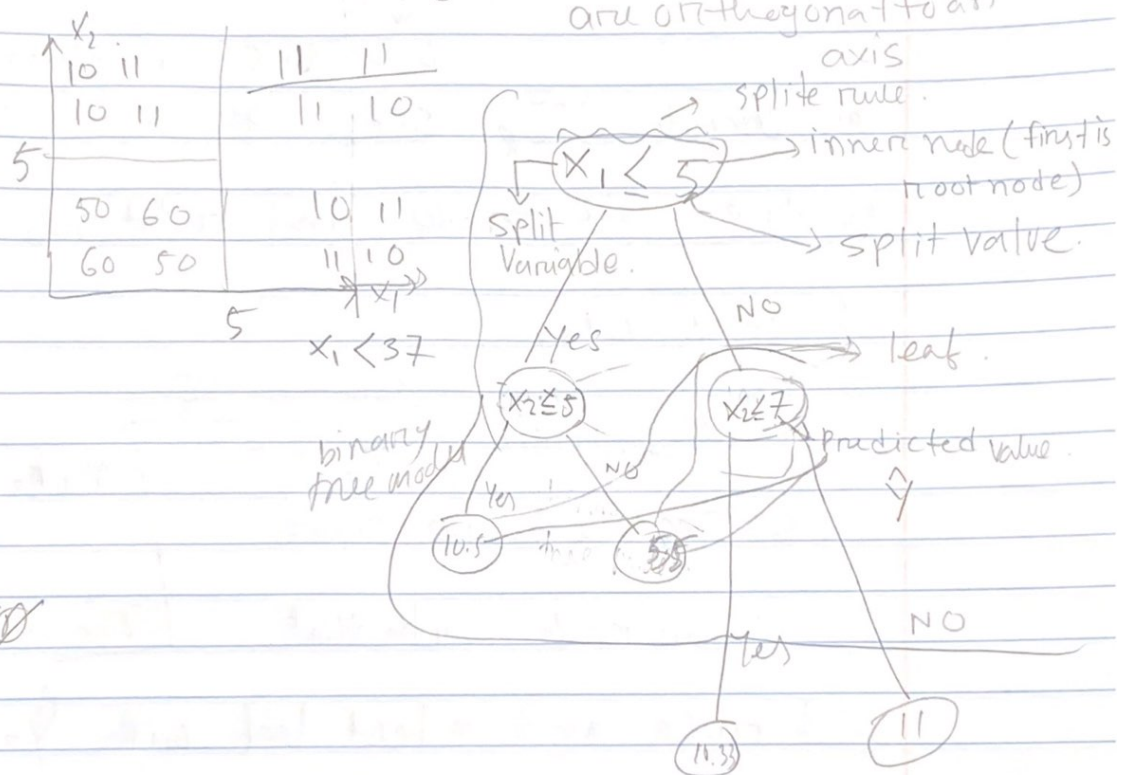opt. Model

Select bin size using the OOSSE curve

Next notebook.

~~Leafa~~

In two dimensions $(P=2)$, the bins are squares.

$B = \#$ bins/dim $\Rightarrow B_{tot} = B^2$.

In P dimensions $B_{tot} = B^P$ which $> n$ very fast

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



$\rightarrow$ split rule.

$\rightarrow$ inner node (first is root node)

$\rightarrow$ split value.

$\rightarrow$ leaf.

Split Variable.

$x_1 < 37$

binary tree model

predicted value

X

= Union of mutually exclusive,

Collectively exclusive, Possibly infinitely

large hyper____

Regression Tree algorithm,

⓪ Let dataset be all data

① Consider every possible orthogonal-to-axis split $X_j \leq X_{(i)j}$ . $j = 1 \cdots p$, $i \in 1, \cdots n-1$.

$\downarrow$

ordered/sorted values.

and compute $SSE_\ell$, $SSE_r$ the SSE's in the putative left mode and right node. Select the rule where,

$$SSE_{weighted} = \frac{n_\ell SSE_\ell + n_r SSE_r}{n_\ell + n_r}$$

is smallest. ie create

inner node with that

$n_\ell$ = # observation in left set

$n_r$ = # " right "

split rule and a left leaf with $\hat{y} = \bar{Y}_\ell$ and a right leaf with $\hat{y} = \bar{Y}_r$.

② if $n_\ell > N_0$ then set dataset = left

and run step 1.

if $n_r > N_0$ " " = right "

$N_0$ is a hyper meter. If $N_0$ is small e.g

$N_0 > 1 \Rightarrow g$ is overfit

If $N_0$ is large $\Rightarrow$ Underfit model. How to Pick $N_0$? Use 3-fold selection.