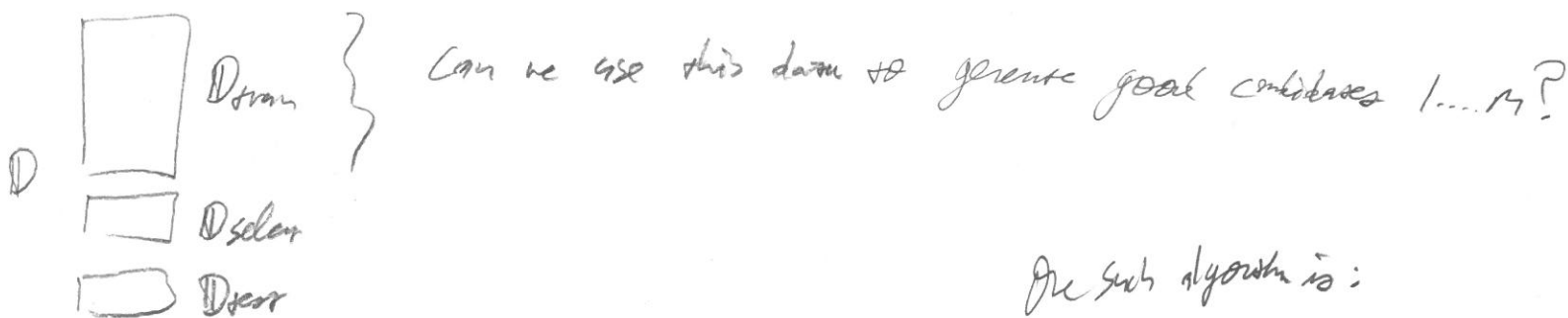


Premisely, we consider "model selection" e.g

Mod 1: $g_1 = b_0 + b_1 x$ from $A_1, \mathcal{H}_1 = \{ \dots \}$
 Mod 2: $g_2 = b_0 + b_1 x + b_2 x^2$ from $A_2, \mathcal{H}_2 = \{ \dots \}$
 Mod 3: $g_3 = b_0 + b_1 x + b_2 x^2 + b_3 \log(x)$ from $A_3, \mathcal{H}_3 = \{ \dots \}$
 :
 Mod M: $g_M = b_0 + b_1 x + b_2 x^2 + b_3 x^3$ from $A_M, \mathcal{H}_M = \{ \dots \}$

We have discussed where models 1...M come from. Above - they
 were made up. Is there something smarter? This is...



The such algorithm is:

For linear models, this is usually done iteratively. You start with
 a huge # of candidate predictors (e.g polynomials, logs, interactions, etc.),
 then iteratively add the "best one". (Difficult definition of "best")

but if you add too many... you overfit!!!

SO... use D_select to stop once you see yourself overfit



Just threshold from
 Error \rightarrow over
 fit

How
 much
 backwards
 stepwise
 work?

This is called "forward stepwise" linear modeling. DEMO

The main problems

- * ① You still need an intelligent selection of predictors to choose from at each iteration → are you sure you can specify this?
- ② The model will still be linear. Although you can alter the A to do non-linear models.

if
linear!

"parametric regression" if $y \in \mathbb{R}$
"parametric distribution" $y \in \{0,1\}$

How about something radically new? When we specified $\mathcal{H} = \{\vec{w} \cdot \vec{x} : \vec{w} \in \mathbb{R}^{n+1}\}$ we made a "parametric assumption" i.e. the model will be of this form which is fixed. Would it be nice if \mathcal{H} can adjust flexibly and allow for more complexity if it increases? This is called "non-parametric regression". Here, the model does not take a prespecified form (like $w_0 + w_1 x$) but is constructed w.r.t. to information in the data \Rightarrow the data provides the model space & the model estimates.

Many creative ideas for this! Let's discuss one such idea

Classification & Regression Trees (CART, 1984)

- Classification trees $y \in \{1, \dots, K\}$
- Regression trees $y \in \mathbb{R}$

Let's talk about regression trees

WAIT do
illustration first

Recall from midsem 2... $\sin(x)$, but no ^{linear} conclusion

3



we see $f(x)$. It looks like a sine curve.

If $\mathcal{H}_1 = \{w_0 + w_1 \sin(w_2 x) : \vec{w} \in \mathbb{R}^3\}$ we would likely do very well with 3 parameters
 A? we can still use LS. but we may need to use a numerical optimizer

If $\mathcal{H}_2 = \{w_0 + w_1 x + w_2 x^2 : \vec{w} \in \mathbb{R}^3\}$ we would likely do very poorly with 3 params.

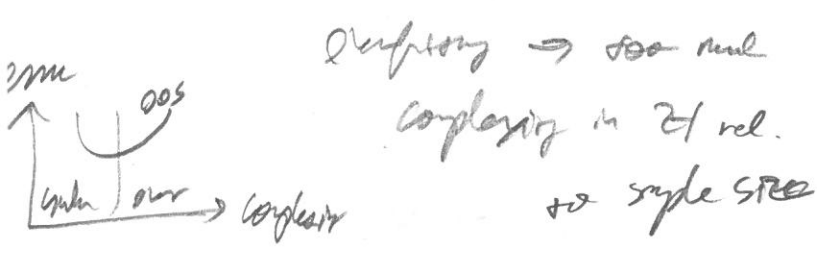
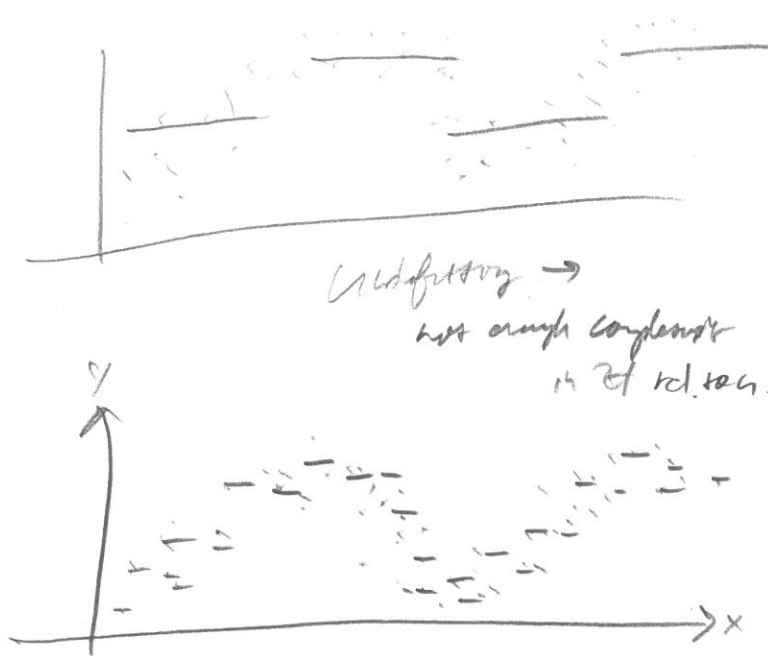
In one dimension this is easy to pick \mathcal{H}_1 over \mathcal{H}_2 . In multiple dimensions, VERY difficult!!

What if we had an algorithm that did the following:



So $g(x)$ is a collection of constant subfunctions where each has a small interval. The union of all intervals = X .

You can make even smaller domains for better resolution.



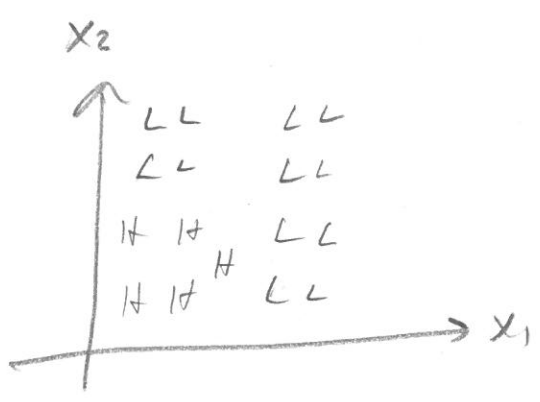
ex: large models
What ^{type of} error is produced here?

misspecification error
(it not express any)

ex: small models
What type of error is produced here?

learning error
(not enough data in each
single model so make a good
decision about where to put
the constant value)

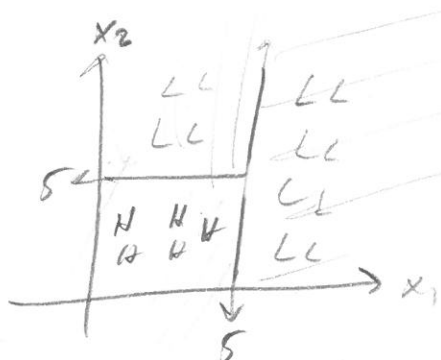
We will discuss how to find a nice balance between underfitting/overfitting later. First let's talk about what happens in $p=2$



let $y = R$
but for purposes of illustration,
let's assume L : low value
 H : high value

What do models look like in 2-dim? Shapes!

Let's make a simplification: shapes must be rectangles whose sides are parallel to the axes. What "model space" makes most sense here?



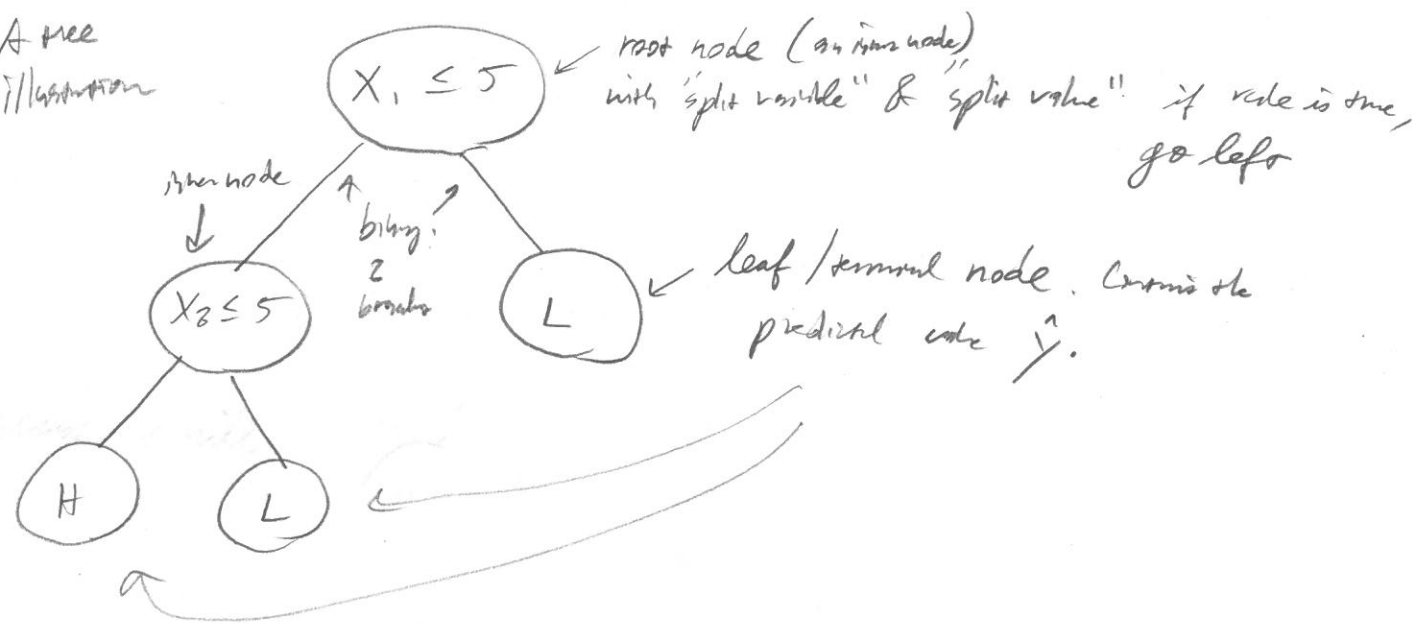
Who is this model? 3 regions (all labels in size)

$$g(\vec{x}) = L \mathbb{1}_{x_1 > 5} + L \mathbb{1}_{x_1 \leq 5} \mathbb{1}_{x_2 > 5} + H \mathbb{1}_{x_1 \leq 5} \mathbb{1}_{x_2 \leq 5}$$

It's actually a linear model...

Is there a way to visualize this? A "binary tree"

A tree illustration



Is all \mathcal{X} covered for? Yes. There is a $\hat{y} \forall \vec{x} \in \mathcal{X} = \mathbb{R}^2$.

Where is the rectangle - parallel-to-axis restriction seen?

Rules must be of the form " $x_j \leq m$ " creating two rectangles parallel to the x_j direction.

$\mathcal{H} = \{ ? \}$ good question ... the...

What is \mathcal{A} ? There are ∞ possible models of MANY!!

this form. How to select (1) splits (2) predictions. One such \mathcal{A} ?

model structure \uparrow specific \uparrow param estimates.

① Begin with all ^{train} data X, Y

② For every possible split at the current node divide the

data into X_L, Y_L & X_R, Y_R and calculate

$$SSE_L = \sum (y_L - \bar{y}_L)^2, \quad SSE_R = \sum (y_R - \bar{y}_R)^2$$

↑
over # of
data pts under the split

③ Find the split with the lowest overall avg SSE

$$SSE_{avg} = \frac{n_L SSE_L + n_R SSE_R}{n_L + n_R}$$

④ Create the split, split data into two nodes

⑤ Repeat steps 2-4 until "STOP"

STOP: node has $\leq N_0$ data pts inside.

There are many, many variations of the above.

If $N_0 = 1 \dots$ tree is grown to fit a separate parameter for each data pt. $R^2 = ?$ 100%! Then tree is "pruned" back to not overfit.

How is N_0 picked? Model selection procedure

Kind of like backwards stepwise regression.

