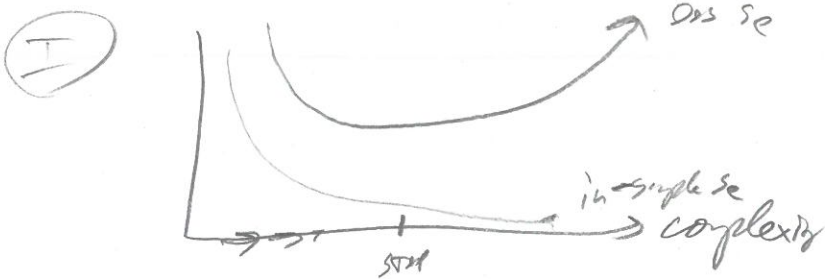


What is a common use of this system?



Use this to increase complexity unit-by-unit until Se is Δ then increases.

This is called ^{greedy} stepwise model selection. This is how it

works: ① pick a set of p linear predictors: $x_1, \dots, x_p, x_1^2, \dots, x_p^2, \dots, h(x_1), \dots, h(x_p), x_1 x_2, x_1 x_3, \dots, x_{p-1} x_p$, etc.

① For each of the p predictors, try each one by adding it to the model and select the model with maximum reduction in in-sample Se .

② For each of the remaining $p-1$ predictors, return to step 1. Sweep of Se on validation set increases.

③ Using these p features, build model on $D_{train} \cup D_{validation}$ and test on D_{test} to get an out-of-sample estimate.

④ \dots build model on D and step g .

(II) Select Hyperparameters. ✓ do this first

Recall SVM which minimized:

$$AHE + \lambda \|\tilde{w}\|^2$$

How to pick λ ?

[2]

Hypothesis selection aka "tuning" which is same as model selection

- ① Select grid of λ 's. $\lambda_{grid} = \{ \overset{.003, .03, .3, 3, 33}{.00, .01, .1, 1, 10, 100} \}$
- ② For each λ is a different A . Fit q models.
- ③ Select model with smallest error in validation set.

Problem 2 ↑

Previously we considered $\{g_1, \dots, g_m\}$ and used cross-validation to select a good model, ^{for min. in-sample error.} This grid method does cross-validation to obtain an honest estimate. But how to pick the models $\{g_1, \dots, g_m\}$?

How about we create an A that flexibly makes \mathcal{A} more complicated? so more $n \rightarrow$ more d.f.. And this A is immune to overfitting. This is the goal of non-parametric machine learning. We will study tree models, specifically

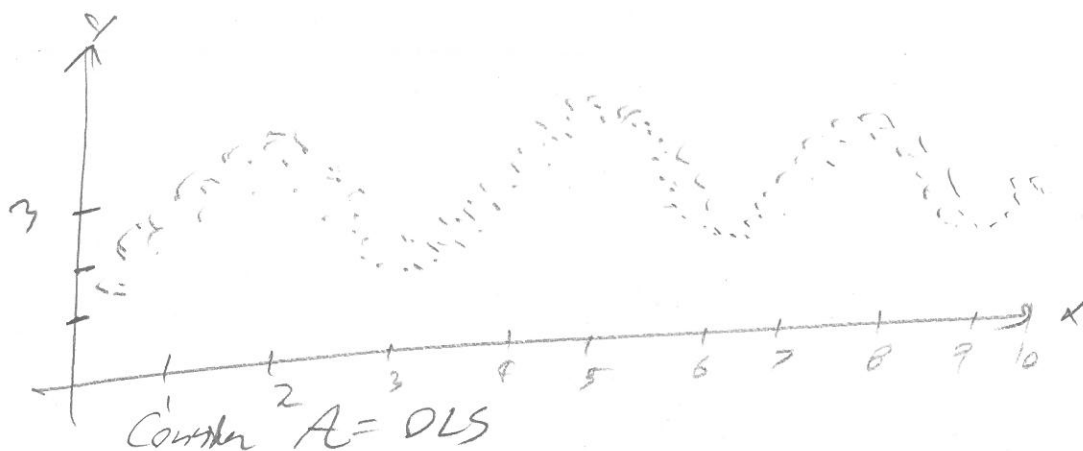
Classification & Regression Trees (CART, 1984)

$$\downarrow$$
$$y = \{1, \dots, L\}$$

$$\downarrow$$
$$y \in \mathbb{R}$$

Imagine done as:

LS



$$\mathcal{H}_1 = \{w_0 + w_1 x : \vec{w} \in \mathbb{R}^2\}$$

$$\mathcal{H}_2 = \{w_0 + w_1 x + w_2 x^2 : \vec{w} \in \mathbb{R}^3\}$$

$$\mathcal{H}_3 = \{w_0 + w_1 \sin(w_2 x) : \vec{w} \in \mathbb{R}^3\}$$

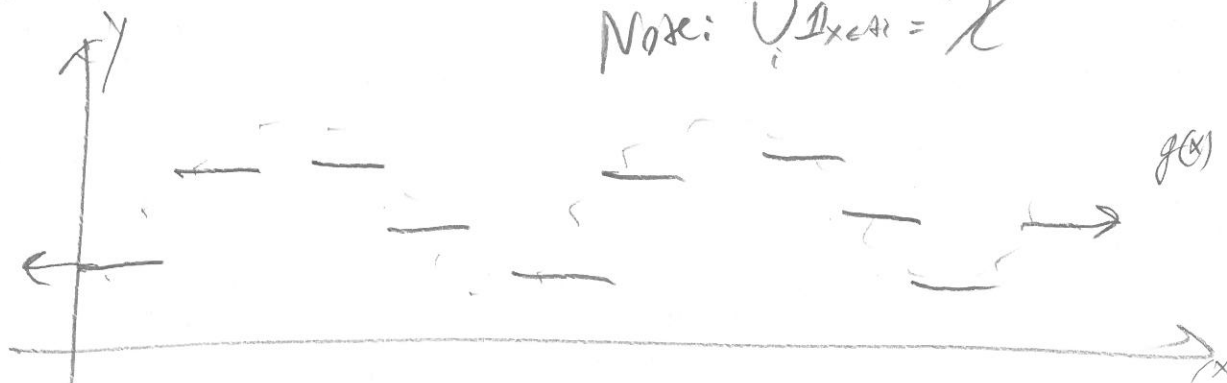
Easy to see \mathcal{H}_3 would win! But in high p - not easy!!

Now consider

$$\mathcal{H}_f = \{w_1 \mathbb{1}_{x \in (0,1)} + w_2 \mathbb{1}_{x \in (1,2)} + \dots + w_{10} \mathbb{1}_{x \in (9,10)} : \vec{w} \in \mathbb{R}^{10}\}$$

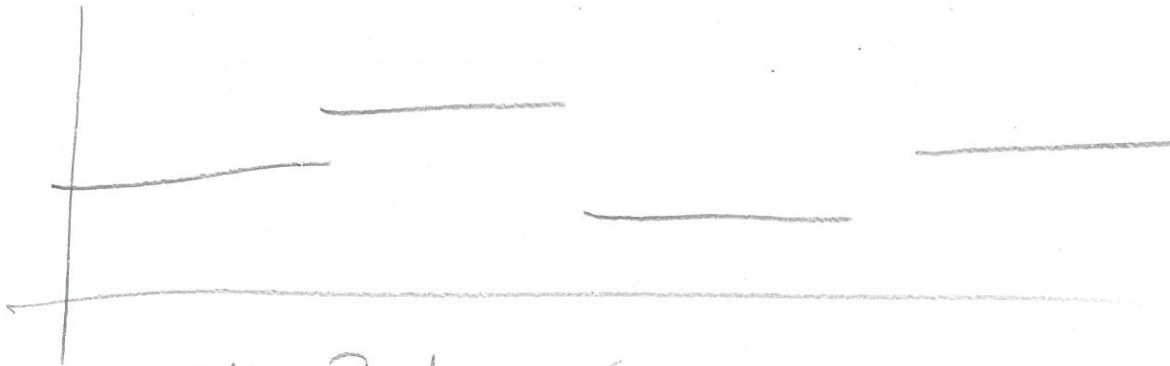
the fit b_1, \dots, b_{10} is given below. The estimates are just \bar{y} within the bins.

Note: $\bigcup_i \mathbb{1}_{x \in A_i} = \mathcal{X}$



Let's make less ^{slow} parameters by using large bins

19



Underfitting? Yes (high bias error)

Max ^{slow} parameters.



Overfitting? Yes (extreme error) Not enough data in each bin to make a good decision

How to pick bin size? Use the model selection routine just discussed.

