

Recall

$$y = g(x) + \underbrace{h(x) - g(x)}_{\text{estimation error}} + \underbrace{h(x) - f(x)}_{\text{misspecification error}} + \underbrace{\epsilon(x) - f(x)}_{\text{error due to randomness}}$$

$\xi(x)$

$e(x)$

Previously, we measured $e_i = y_i - \hat{y}_i = y_i - g(\hat{x}_i)$

for $i \in \{1, \dots, n\}$ since D has n observations inside. These e_1, \dots, e_n are "in-sample residuals". Computations of $\text{SSE}, \text{RMSE}, R^2$ are

all in-sample. We discovered last time that in-sample error metrics can be made arbitrarily more perfect simply by adding additional variables (i.e. raise p).

This is not only the case with $\mathcal{R} = \text{OLS}$, this is a problem for most supervised learning alg's. You can see it more starkly with OLS!

Thus e_1, \dots, e_n are not good estimates of the e here. What is the e there? Out of sample.



We forget to reason... we need to...

Assume t, \vec{z} are the same across time and that the \vec{z} 's are measured the same way. This is called "Stationarity" of the phenomenon. We need this otherwise... prediction won't work!

What are good estimates of e, the model's "error"?

$e_x = y_x - \hat{y}_x$ for an $\vec{z}_x \in D_{\text{tr}}$. These are in-sample estimates?

Why? They cannot be overfit. Overfitting is tendency to memorizing the answers on the ^{practice} test. "Fitting" is learning the theory ^{from the old answers} to be able to answer new questions. "Unlearning" is not learning as much as you overfitting. Studying the practice test.

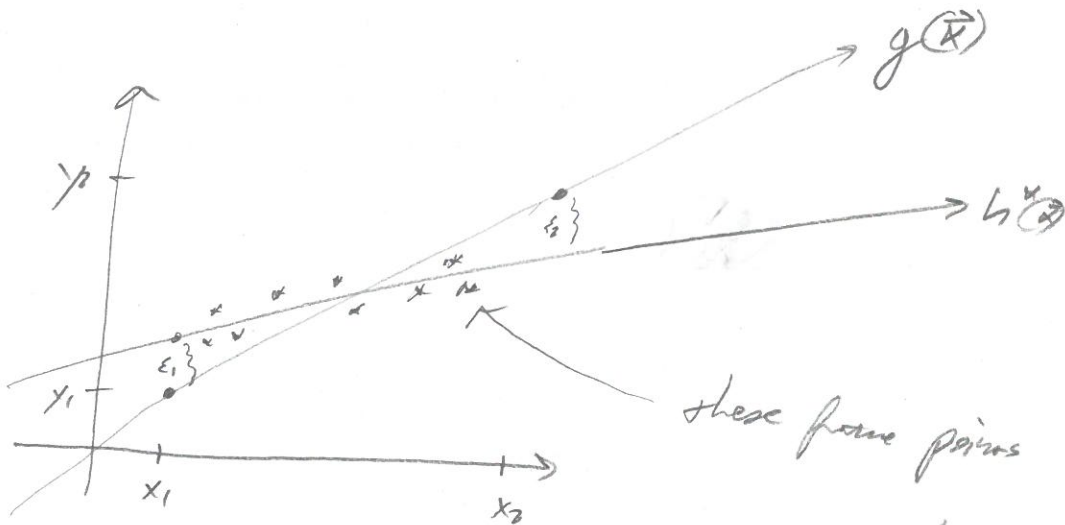
Overfitting... Extract some of the understandable structure is E and attributed it to model structure.

Overfitting in one dimension with intercept...

(3)

Recall

$$y = h^*(x) + \varepsilon, \quad h = 2, \quad p = 1 \Rightarrow \text{perfect fit}$$



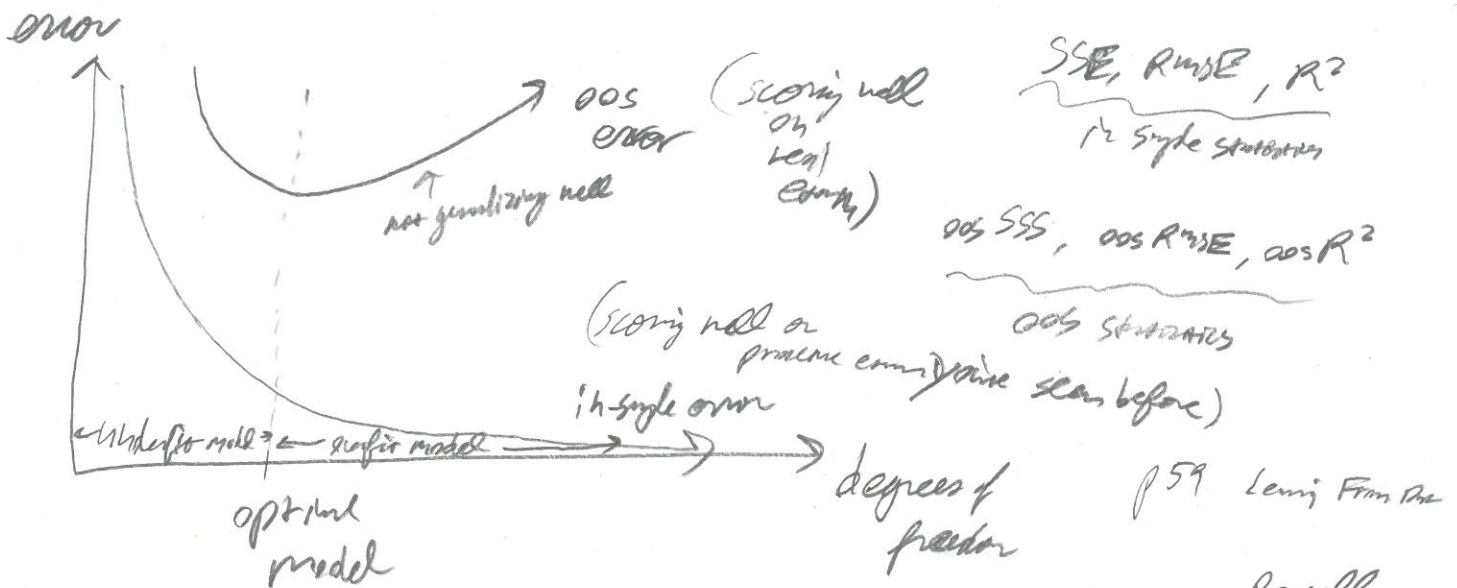
these frame points are further from g than h^* .
do's generalize well.

$$y = g + \underbrace{(h^* - g)} + \underbrace{(g - h^*)} + \underbrace{(\varepsilon - \varepsilon)}$$

these three errors are all generalization errors
since they can be made $\rightarrow 0$ simultaneously

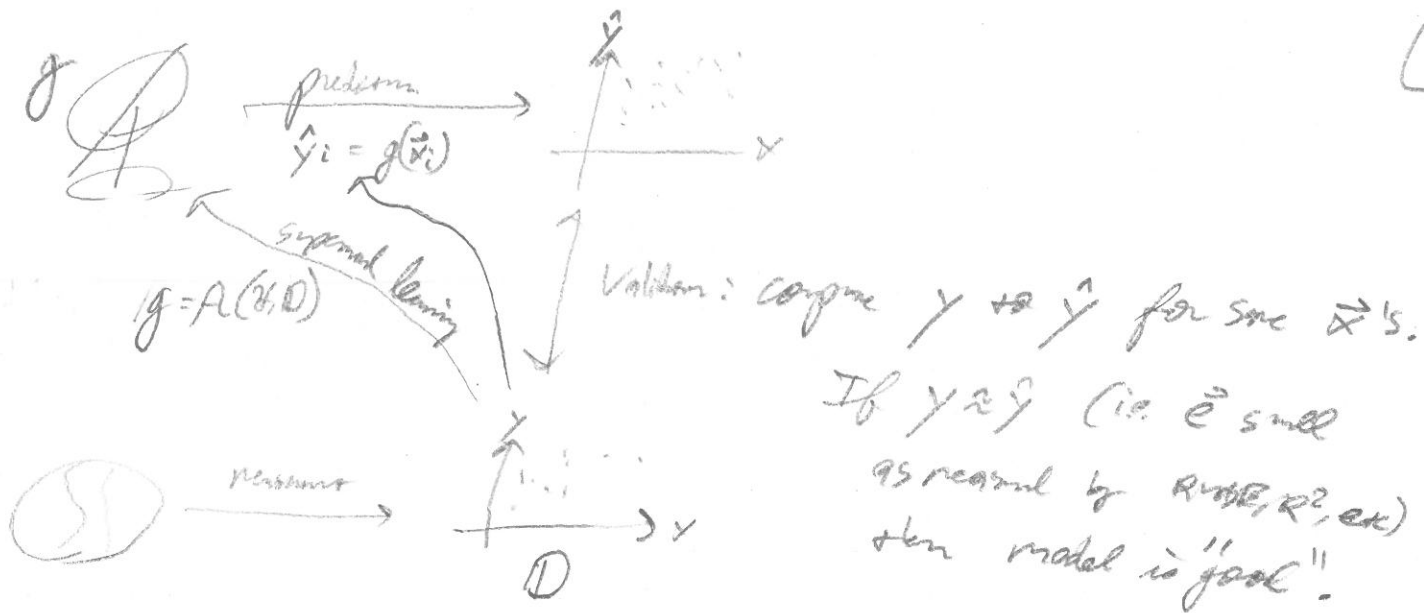
Overfitting does not change h^* or ε or t , only g .

\Rightarrow Overfitting increases estimation error. (Recall)



Recall

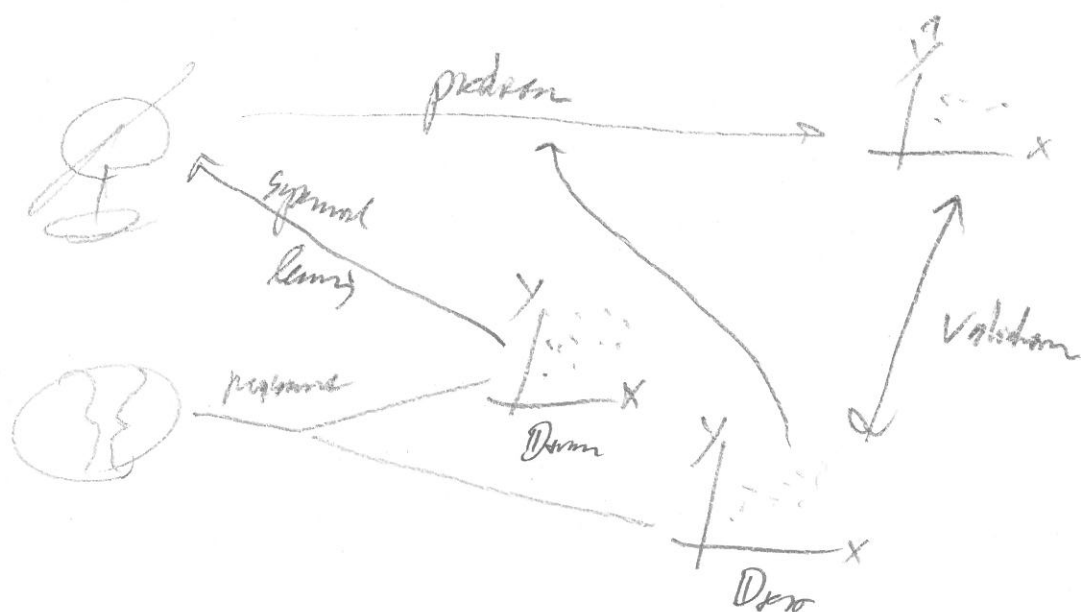
In-sample error $\rightarrow 0$. It's like memorizing the answers! Random y_1, \dots, y_n . KAN $K=1$!



In-sample validation: Use D to build g . Get \hat{y} from D and compare to y . BAD! Can be overfit.

Solution: split $D = D_{\text{train}} \cup D_{\text{test}}$

D_{train} becomes the data used to build the model i.e. $g = A(\vec{x}, D_{\text{train}})$
 and D_{test} becomes the data used to validate

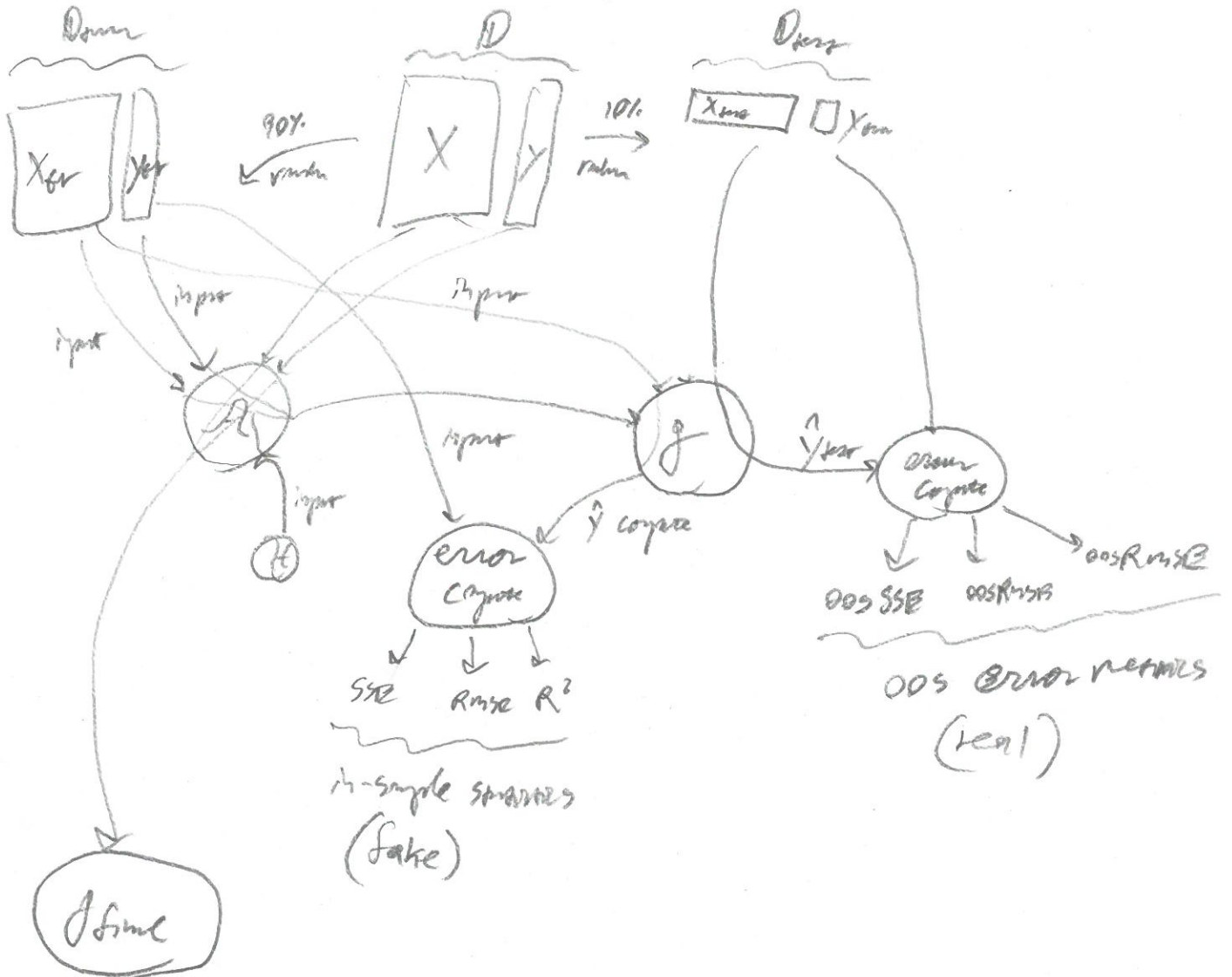


Final model: uses all D .

typical split e.g. proportion left out for testing is 10%, 20%, 5% 5

Usual notation:

$$K = \frac{1}{\epsilon} \quad \text{e.g. } \epsilon = 10\% \Rightarrow K = 10. \text{ Then}$$



the model used for prediction

$f = f_{in}$ No. f_{in} is better on avg than f . So oos error metrics are an upper bound / conservative estimate

Which error does f_{in} have less of? Estimation error! It is 11% higher!

Not huge gain... but why not take the cookie?