

p59 (Learning from Data)

In-sample error: SSR, S_e computed from D_{train}
 Out-of-sample error: \hat{SSR}, \hat{S}_e computed from random D^*
 (OOS)
 your score on practice exam that you get to see
 ' ' ' ' the real exam ' ' ' ' don't ' ' ' '

In-sample error can go to zero since you can merely memorize all $i=1, \dots, n$ answers, the $\{y_1, \dots, y_n\}$. If you're merely memorized, then you haven't learned the theory/concepts/principles which is $f(x)$ or at the very least, $h^*(x)$.
 (18N, 3mls?), ...

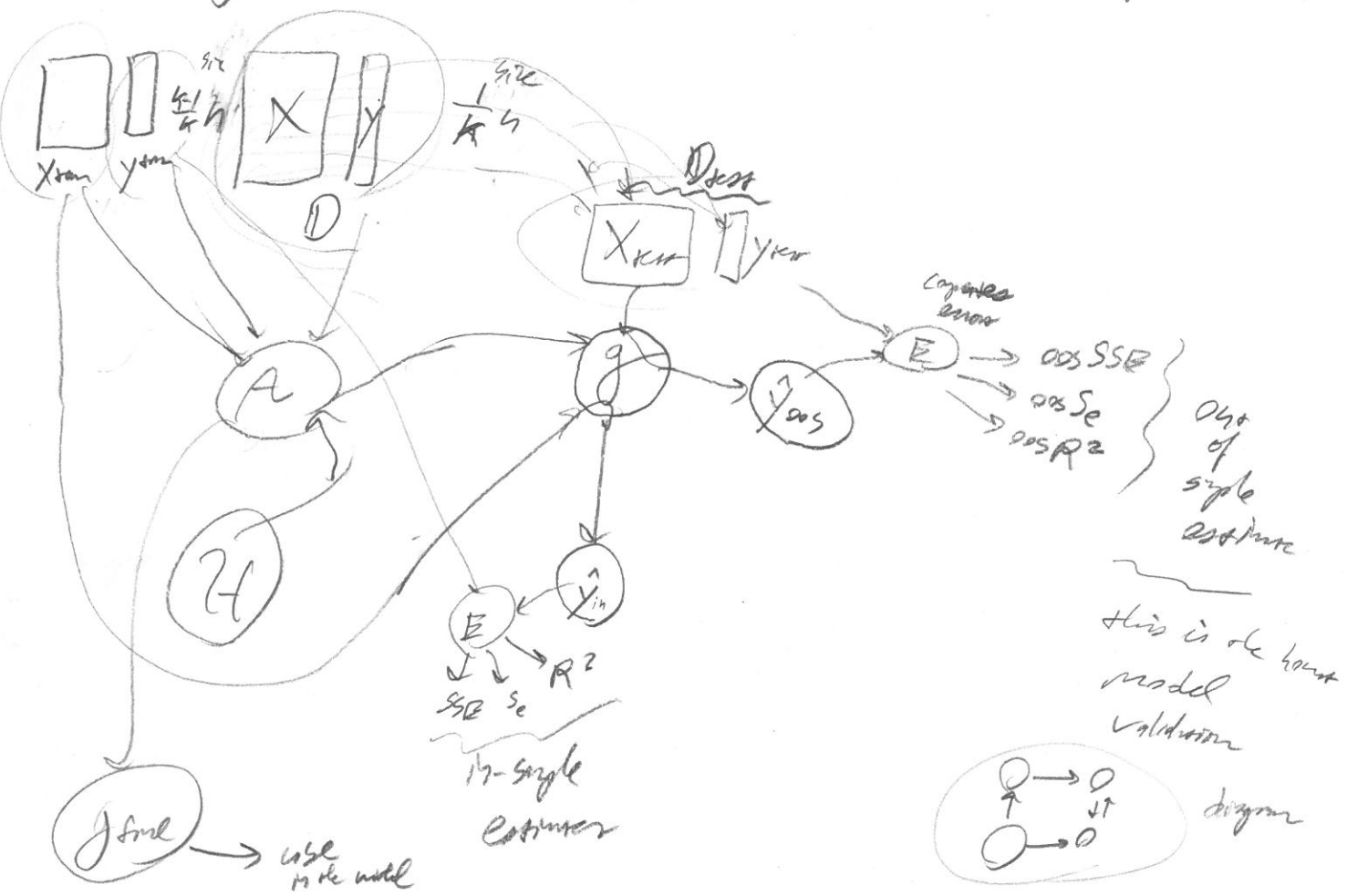
Imagine if Newton only memorized Force-mass-accel data? Instead of memorizing, the law $F=ma$ is the dimension reduction, the principle that is learned.

In-sample error metrics are not honest. How do we get Validation / generalization?
 ED how seen by $A \Rightarrow$ cannot have been incorporated into model f
 Need new data where the y_i 's are known so we can compare to \hat{y}_i 's from f . How do we do this if D is all we got ???

Split $D = D_{train} \cup D_{test}$ randomly.

Typical splits are 90% train & 10% test or 80% train & 20% test.
 $K=10$ $K=5$

Use $g = A(D_{train}, H)$ and use D_{test} to predict oos.

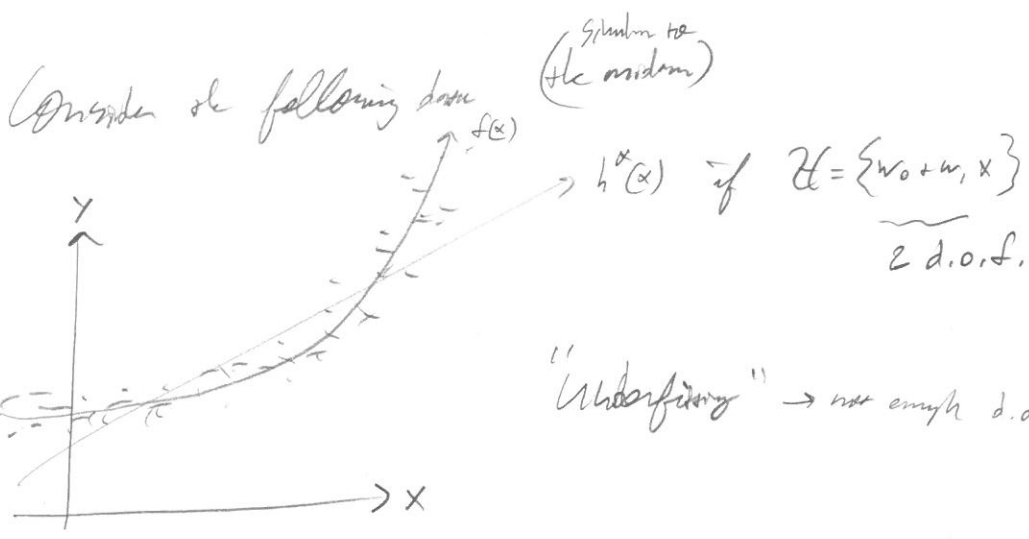


In order for the OOS error estimator to be honest D_{test} can not be looked at until g is fully constructed. And it can be used once to get OOS estimates and can never be looked at again.

It's like passing to the Turing test... once passed... it is never passed again.

For the final fitted model, use all D . This final model is not validated directly but is assumed to be more accurate than g since now it has less estimation error.

DEMO



"Underfitting" \rightarrow not enough d.o.f.

However, there is misspecification error. We know (because we can see) $f(x)$ - this is possible only in one dimension. As a result why does $y \neq f(x)$? $f \rightarrow$ ignore error.
How do we allow $f(x) \in \mathcal{H}$? Make \mathcal{H} "richer". How?
We need to add more terms beyond the linear terms.

Consider the Weierstrass Approx. Thm:

For every cont. function x in region $\in \mathcal{X} = [a, b]$, \exists a polynomial function $p(x)$ s.t. $\forall \epsilon \exists \forall x \in \mathcal{X} |f(x) - p(x)| < \epsilon$.

This means any cont. function can be approximated with a polynomial.
Let's begin with polynomials of degree 2. $\Rightarrow \mathcal{H} = \{w_0 + w_1 x + w_2 x^2 \mid w \in \mathbb{R}^3\}$
In the above example $p=1$

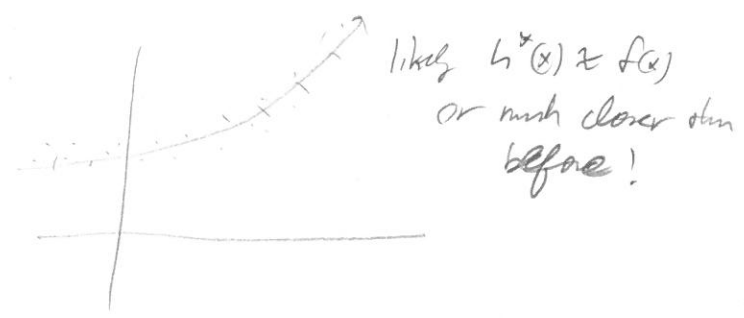
Since \Rightarrow add more d.o.f. \Rightarrow add poly degree two d.o.f.

$$X = \begin{bmatrix} 1 & x_{11} \\ \vdots & \vdots \\ 1 & x_{1n} \end{bmatrix} \Rightarrow X = \begin{bmatrix} 1 & x_{11} & x_{11}^2 \\ \vdots & \vdots & \vdots \\ 1 & x_{1n} & x_{1n}^2 \end{bmatrix}$$

$p=1$ $p+1=2$ & full rank $2p+1=3$ & full rank \rightarrow why?

Now use LS as usual! This is why I call this a "linear non-linear model".

$$\vec{b} = (X^T X)^{-1} X^T \vec{y}$$



Can we get better?

Why not approximate f better by using degree 3?

$$X = \begin{bmatrix} 1 & x_{11} & x_{11}^2 & x_{11}^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{1n} & x_{1n}^2 & x_{1n}^3 \end{bmatrix}$$

What could happen?

DEMO

\Rightarrow If n is large, not really... the x^3 problem would function as one nonsense predictor... not so bad.

What about many more polynomial terms? \Rightarrow BAD

he just saw $n=5$ and a 4-degree polynomial fits it. What?

Vandermonde matrix

$$\underbrace{\begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 & x_1^4 \\ 1 & x_2 & x_2^2 & x_2^3 & x_2^4 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_5 & x_5^2 & x_5^3 & x_5^4 \end{bmatrix}}_X \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_4 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_5 \end{bmatrix}$$

Is this invertible? $\det(X) = \prod_{i=1}^n \prod_{j=i+1}^n (x_j - x_i) \neq 0$ only if all x_i 's unique

Yes n data pts can be fit by a $n-1$ degree polynomial.

Instead of polynomials, could you make logs? Sines/cosines?

Exponential? Others? Yes! The sky is the limit! Hopefully we

can talk more about alternatives to polynomials...