

M342W

3/10

$$y_i = g(\vec{x}) + \underbrace{h^*(\vec{x}) - g(\vec{x}) + h^*(\vec{x}) - f(\vec{x}) + t(\vec{z}) - f(\vec{x})}_{e(\vec{x})}$$

errors:                      estimation                      misspecification                      ignorance

$$e_i = y_i - \hat{y}_i = y_i - g(\vec{x}_i) \quad i = 1 \dots n, \text{ rows of } \mathbb{D}$$

These  $e_1, \dots, e_n$  are called "in-sample residuals" because they come from  $\mathbb{D}$  (the sample). Thus, SSE, RMSE,  $R^2$ , SSR are all "in-sample".

Thus, they can all be "fake". Why? Because we saw that if  $p$  goes up, i.e. we make up new  $\vec{x}$ 's, we can get any answer we want.

We can make  $R^2$  arbitrarily high and RMSE arbitrarily low. Thus,  $e_1, \dots, e_n$  are not honest estimates of our prediction errors/prediction accuracy.

Imagine  $\begin{array}{c} \mathbb{D} \text{ past} \\ \mathbb{D}_* \text{ future} \end{array} \downarrow \text{time}$

A good way to honestly estimate prediction errors/prediction accuracy is to compute  $\vec{e}_* = \vec{y}_* - \vec{\hat{y}}$  i.e. on data you did not use to build your model( $g$ ).  
 $\vec{\hat{y}}_*$  cannot be overfit.

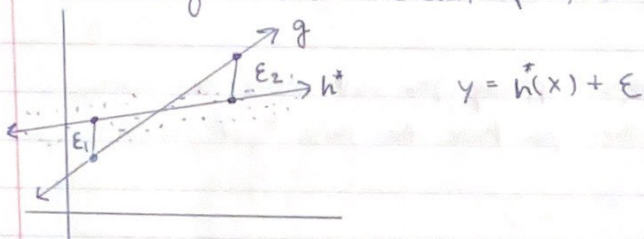
For this to be true, we need another assumption, "stationarity" i.e. that  $t(\vec{z})$  stays the same and the relationships between the  $\vec{z}$ 's and the  $\vec{x}$ 's remain the same which means the function  $f$  doesn't change with time  $t$ . An example of a non-stationary relationship is stock prices ~~vs.~~ explained by some variable  $\vec{z}$  that are causal at one moment.

What is underfitting? Not learning as much as you can with the information you're given?

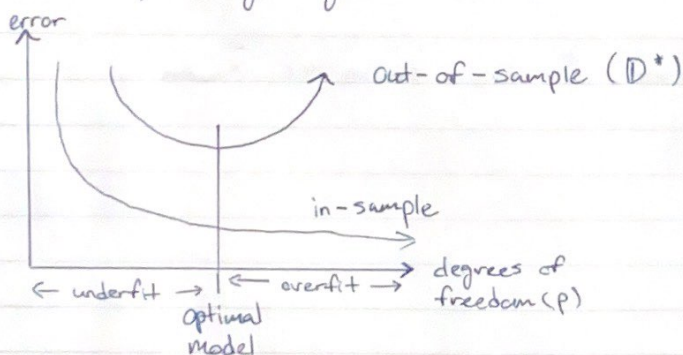
Overfitting - estimation error

Underfitting - misspecification error

Overfitting in one dimension ( $p=1$ ) with two observations ( $n=2$ ).



Overfitting is fitting the  $\epsilon$ 's which you know is a bad idea since  $\epsilon$  = misspecification error + error due to ignorance. Fitting any degree of  $\epsilon$  leads to a poorer model. Overfitting doesn't change  $h^*$  or  $f$ , it only changes  $g \Rightarrow$  it's error in estimation.



We need out of sample (honest) metrics:

$$\text{out RMSE} = \sqrt{\frac{1}{n_* - (p+1)} \text{cosSSE}} \quad \text{or} \quad = \sqrt{\frac{1}{n_*} \text{cosSSE}}$$

$$\text{cosSSE} = \sum_{i=1}^{n_*} e_{*i}^2$$

$$\text{cos } R^2 = 1 - \frac{\text{cosSSE}}{\text{SST}_*} = \frac{\sum (y_{*i} - \bar{y}_*)^2}{\text{SST}_*}$$



oos error metrics require computation on data from the future (which you don't have). So how can we possibly compute oos error metrics? So... assume stationarity and partition our original data set into:

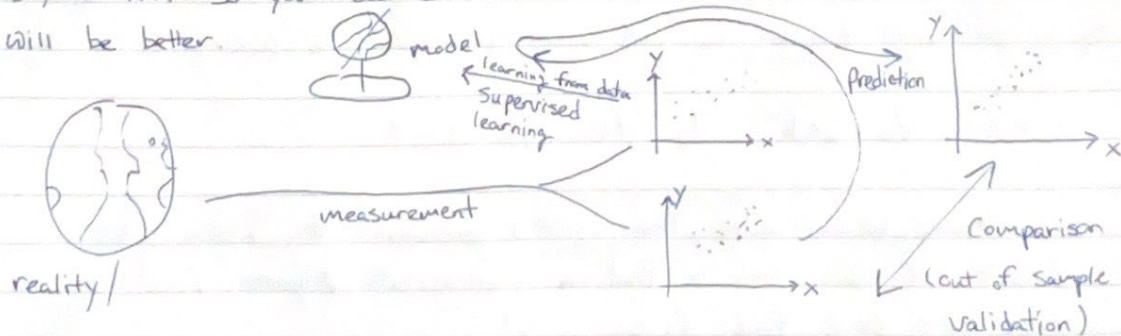
$$D = D_{\text{train}} \cup D_{\text{test}}$$

Then use  $D_{\text{train}}$  the way we've been using  $D$  all along which is use it to create your model:  $g = A(D_{\text{train}}, \{L\})$ , and then use  $D_{\text{test}}$  to compute the honest oos error metrics to inform us of our model's predictive accuracy in the future (i.e. give us an estimate of our generalization error). Two issues:

- 1) What proportion of the data do you use for training?
- 2) how to do the split?

- 1) really has no answer but people generally use 80% or 90%.
- 2) usually done randomly in case there's a time trend in  $D$ .

$n_{\text{train}} < n$ . is there a cost to this? Yes, estimation error increases and thus, the generalization error estimate (oos error metrics) will look a tad worse than they will be when using your models (since the model you will use in the future will be built with all  $D$ , all  $n$ ). So your oos error is conservative; the real oos error will be better.



$$h^*(\vec{x}) - g(\vec{x}) = \vec{x}\vec{\beta} - \vec{x}\vec{b} = \vec{x}(\vec{\beta} - \vec{b})$$

Estimation error in OLS. As  $\vec{b}$  diverges from the  $\vec{\beta}$ , the estimation error goes up. One way to measure it is via  $\|\vec{\beta} - \vec{b}\|^2$

The usual notation is to specify a  $K$  where  $1/K :=$  proportion in the test set. So default  $K = 5$  or  $10$

$$K = 5 \Rightarrow 1/5 = 20\% \text{ Dtest}$$

In the above procedure, there is only one oos validation comparison.

If you see that you overfit by seeing that e.g.  $\text{RMSE} \ll \text{oos RMSE}$ , it is too late! You can't correct the model honestly anymore.

So... this is a big problem and we need to solve it.

But first let's talk about reducing misspecification error.