

# lecture 16



$$\frac{1}{k} = \frac{n_{test}}{n}$$

$$\Rightarrow k = \frac{n}{n_{test}}$$

$D$  only a sample of the full distribution  $\langle X, Y \rangle$ .

$$n_{test} \in \left\{ 1, 2, \dots, \frac{n}{3}, \frac{n}{3}, \dots, \frac{n}{5}, \frac{n}{10}, n-1 \right\}$$

$$\Rightarrow k \in \left\{ n, \frac{n}{2}, \dots, 2, 3, 5, 10, \frac{n}{n-1} \right\}$$

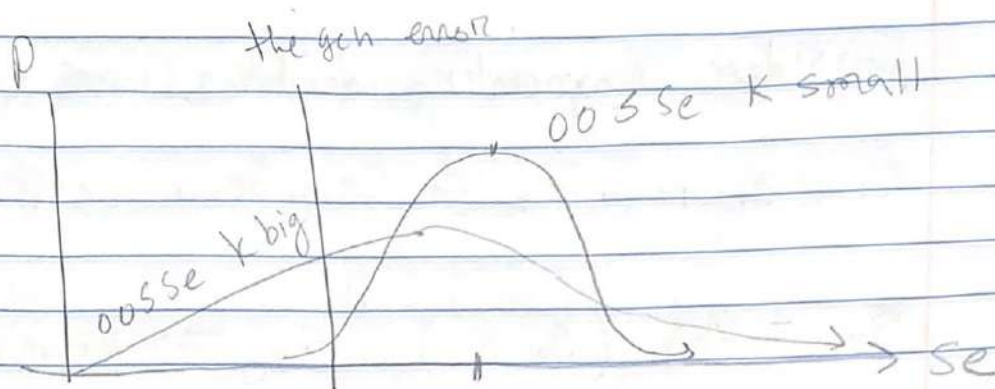
most common choices for  $k$ .

What's the trade off?

We want Validation! We use  $D_{test}$  to estimate "generalization error".

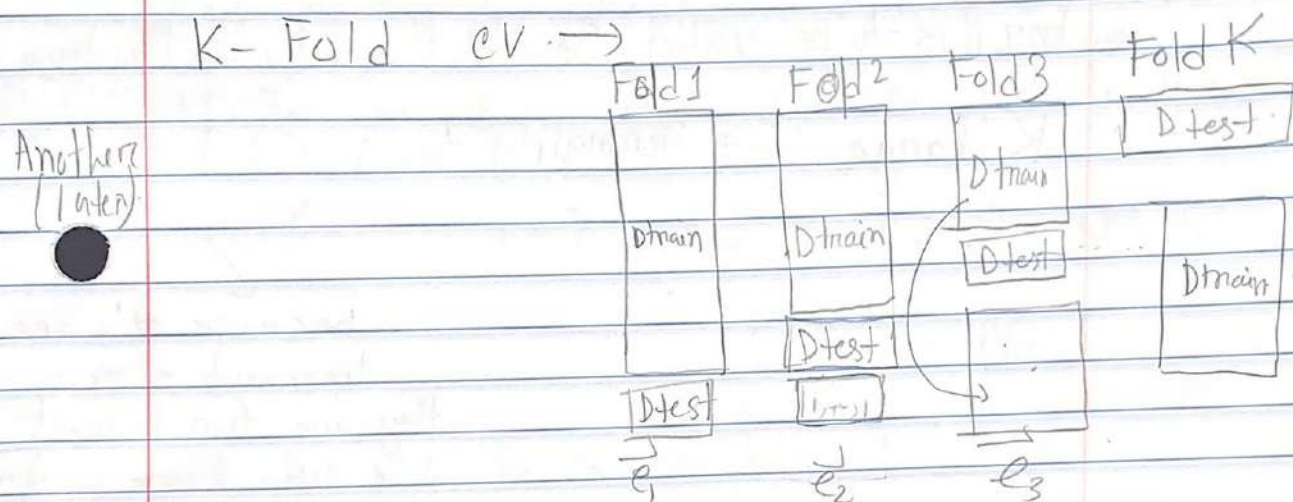
AKA model performance in the future. If  $n_{test}$  is small, the estimate is very variable but

$n_{test}$  is big... my expected estimate is closer to true performance at final



Cross Validation:- do many train-test splits.

K-Fold cv.  $\rightarrow$



Overall  $\parallel$  K folds, each observation is oos once

$$\text{Let } \vec{e}_{cv} = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

has length  $n$ . Now compute oos metrics on  $\vec{e}_{cv}$ .

if  $K = n$  ( $n$  folds) AKA "leave-one-out"

Cross Validation " (LOOCV)



Consider Computing metrics on  
 $\vec{e}_1, \vec{e}_2, \dots, \vec{e}_k \Rightarrow se_1, se_2, \dots, se_k$

$\bar{se} \approx \text{Avg} (se's), S_{se} \approx \text{Stddev} (se's)$

$$CI_{\sigma_0, 95\%} [\bar{se} \pm 2 S_{se}]$$

For this to be valid;  $se \text{ iid } N(,)$ ; Is this true?  
interpretable? NO.

$K \text{ large} \rightarrow \text{normality from CLT.}$

because  $d$ 's are  
dependent since  
they are functions of  
 $g$ 's which are dependent  
via sharing ID traits

Computing <sup>Valid</sup>  $a_n CI$  for generalization error

under general condition is not  
possible.

Given  $D$ , You have many choices of model  
because there are many  $A$ 's and  
 $H$ 's.  $g = \hat{A}(D, \hat{H})$  [change  $A$  and  $H$ ,  
Therefore  $g$  change  $g$ ]

Let's say you have a choice of a  
finite # of models,  $M$ : <sup>How do we</sup> Choose between

$g_1, g_2, \dots, g_m$ ?  $P_{raw} = 1$ .  $A = OLS$   
 $Y = \mathbb{R}$ .

$$g_1 = b_0 + b_1 X.$$

$$g_2 = b_0 + b_1 X + b_2 X^2.$$

$$g_3 = b_0 + b_1 \ln(X)$$

$$g_4 = b_0 + b_1 \mathbb{1}_{X \in [0,1]} + b_2 \mathbb{1}_{X \in [1,2]} + \dots$$

This is the fundamental problem of  
"Model Selection."



In this class, we will provide an idea.

It is not the only idea

\*  $g_1$  in  $f_1$  with  $A_1, H_1$  on  $D_{train}$  and  $Se_1$  is computed on  $D_{test}$ .

\*  $g_2$  is ~~with~~ fit with  $A_2, H_2$  on Dtrain and  $Se_2$  is  
compared D  
test

N gm n n Am, Hm h v v Se

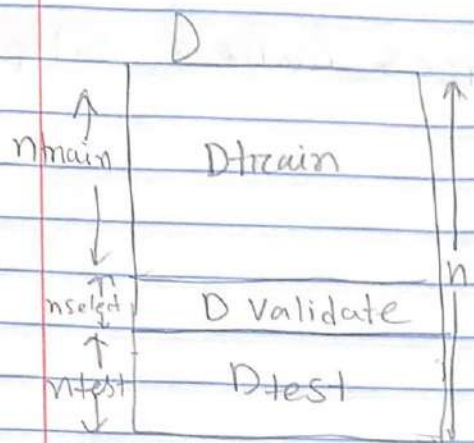
Then select  $g_{m_1}$  which has lowest  $Se$ .

Problems: (1)  $S_e$  has high Variance. So

Pick K reasonable

(2) We have no longer honest Validation.

Se is ~~not~~ honest because it could be overfit since ID test was used many times.



## Procedure

(1) Fit  $g_1$  on  $D_{train}$ ,  
Compute  $Se_1$  on  $D_{select}$

(2) Fit  $g_2$  on  $\dots Se_2$  on  $D_{select}$

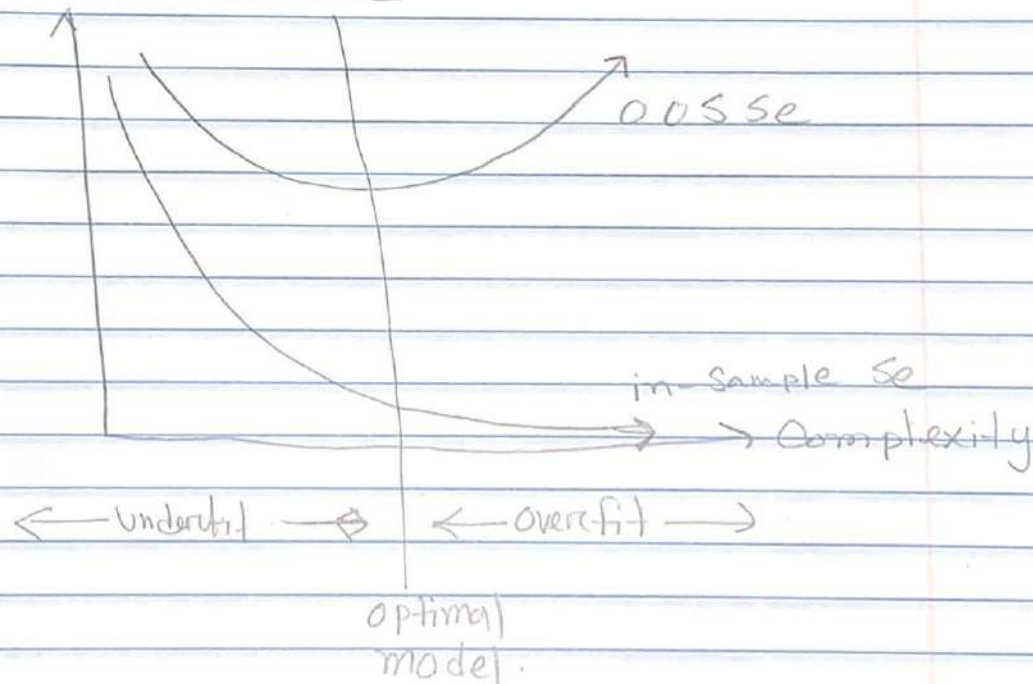
(3) Fit  $g_n$  on  $\dots Se_n$  on  $D_{select}$

(m+1) Final  $m_x = \arg \max \{Se_j\}$

(m+2) Fit  $g_{m_x}$  on  $D_{train} \cup D_{select}$

(m+3) Fit  $g_{final}$  on  $D$

Uses of this procedure:



(1)

Consider  $g_1, g_2, \dots, g_m$  with increasing  
Complexity.

This is Sometimes called Stepwise  
Modeling.