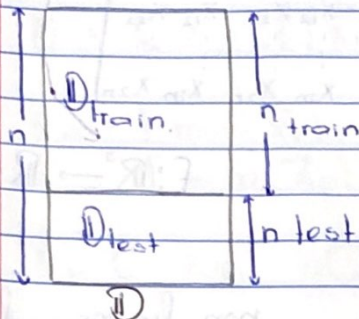


Lecture -16

04/02/2020



$$\frac{1}{k} \text{ train-test} \Rightarrow k = \frac{n}{n_{\text{test}}}$$

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

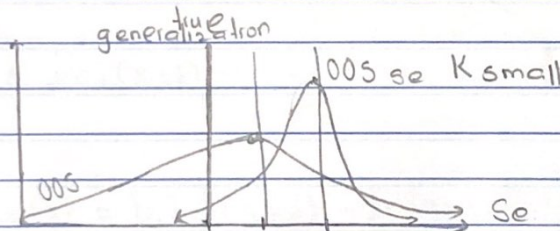
$$\Rightarrow k \in \{n, \frac{n}{2}, \dots, 2, \dots, 3, \dots, 5, \dots, 10, \dots, \frac{n}{n-1}\}$$

Only a sample of the full distribution $\langle \vec{x}, y \rangle$

What's the tradeoff? most common choices for k

We want honest validation, we use D_{test} to estimate "generalization error" AKA model performance in the future.

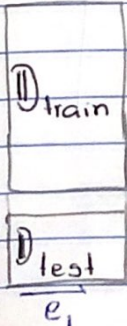
If n_{test} is small, the estimate is very variable but n_{train} is big thus my expected estimate is closer to the true performance of g_{final} .



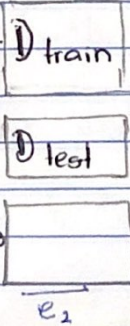
Cross Validation^(cv): do many train-test split.

→ K-Fold cv →
→ Another Center

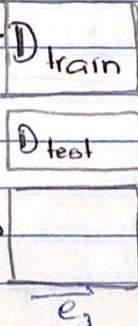
Fold 1



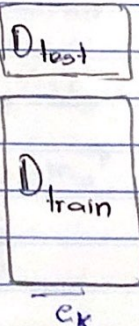
Fold 2



Fold 3



Fold k



Over all k folds, each

observation is OOS once

let $\vec{e}_{cv} = \begin{bmatrix} \vec{e}_1 \\ \vec{e}_2 \\ \vdots \\ \vec{e}_k \end{bmatrix}$ 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 matrices on $\vec{e}_1, \vec{e}_2, \dots, \vec{e}_k \Rightarrow S_{e_1}, S_{e_2}, \dots, S_{e_k}$

$\bar{S}_e = \text{Avg}(S_{e_i}\text{'s}), S_{Se} = \text{Stdev}(S_{e_i}\text{'s}) \quad \text{CI}_{95\%} = [\bar{S}_e \pm 2 S_{Se}]$

For this to be valid, $S_e \stackrel{\text{iid}}{\sim} N(\cdot, \cdot)$. Is this true? And? NO; K -large \rightarrow normality from CLT.

Because \vec{e} 's are dependent. Since they are functions of g 's which are dependant via sharing D_{train} 's.

Computing a valid CI for generalization error under general continuous is NOT possible.

Given ①, you have many choice of model because there are many A 's and many H 's $g = A(D, H)$. Let's say you have choice of a finite # of models, m ; How do we choose between g_1, g_2, \dots, g_m ? e.g. $\text{Praw} = 1$, $A = \text{OLS}$, $y = \mathbb{R}$

$$g_1 = b_0 + b_1 X$$

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

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

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

This is the Fundamental problem of "model selection". In this class, we will provide one idea. It is not the only idea.

g_1 is fit with A_1, H_1 on $\mathbb{D}_{\text{train}}$ and Se_1 is computed on \mathbb{D}_{test}

g_2 is fit with A_2, H_2 on $\mathbb{D}_{\text{train}}$ and Se_2 is computed on \mathbb{D}_{test}

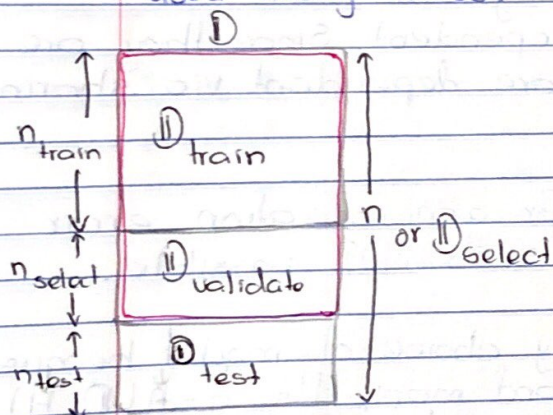
\vdots

g_m is fit with A_m, H_m on $\mathbb{D}_{\text{train}}$ and Se_m is computed on \mathbb{D}_{test}

then select g_{m_*} which has lowest Se .

Problems (1) Se has high variance. So pick K reasonable

(2) We no longer have honest validation. Se is not honest because it could be overfit since \mathbb{D}_{test} was used many times.



Procedure

① Fit g_1 on $\mathbb{D}_{\text{train}}$, compute Se_1 on $\mathbb{D}_{\text{select}}$.

② Fit g_2 on $\mathbb{D}_{\text{train}}$, compute Se_2 on $\mathbb{D}_{\text{select}}$.

\vdots

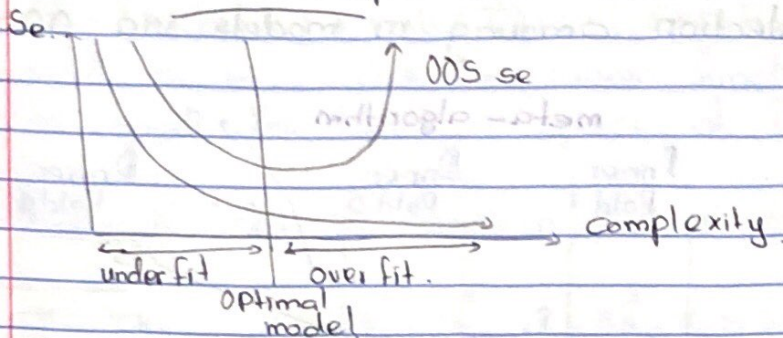
③ Fit g_m on $\mathbb{D}_{\text{train}}$, compute Se_m on $\mathbb{D}_{\text{select}}$.

④ Find $m_* = \text{argmax} \{Se_j\}$

⑤ Fit g_{m_*} on $\mathbb{D}_{\text{train}} \cup \mathbb{D}_{\text{select}}$ and compute Se_{m_*} on \mathbb{D}_{test}

⑥ Fit g_{final} on $\mathbb{D}_{\text{train}}$

Uses of this procedure.



① Consider models g_1, g_2, \dots, g_m with increasing complexity. This is sometimes called stepwise modeling.