

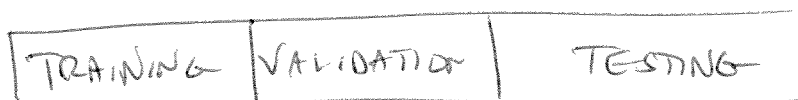
CROSS-VALIDATION

①

So far, we have used model selection criteria that are based on an estimate of the expected gap between $pMSE$ and MSE
e.g. C_p , AIC and BIC.

Now, cross-validation — a direct way of estimating the $pMSE$ itself

Data set split into

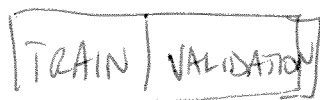


Use TRAINING to fit models M_1, \dots, M_D to data.

Use VALIDATION to compute $pMSE(M_1), \dots, pMSE(M_D)$ and select the model M^* with the smallest $pMSE$

Apply M^* to the TESTING data to obtain a valid estimate of $pMSE(M^*)$

Why not



only? Well, if VALIDATION is used to

find the M^* with the smallest $pMSE$, then this $pMSE$ is a biased estimate of true $pMSE$ ~~error~~ (called selection bias)

It's important to remember that a valid pMSE estimate can only be obtained from an untouched TESTING data.

The pMSE's we get from the VALIDATION data are used to rank the models, not to provide a final estimate of the true predictive performance.

Note, if you only fit one model M_1 to data and evaluate $pMSE(M_1)$ on VALIDATION data, that's OK. It's the ranking of many models that generates the bias.

Now - how to split the data,

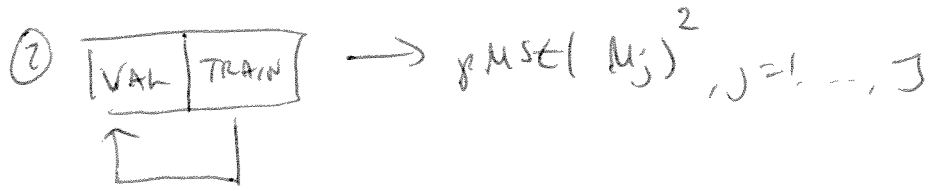
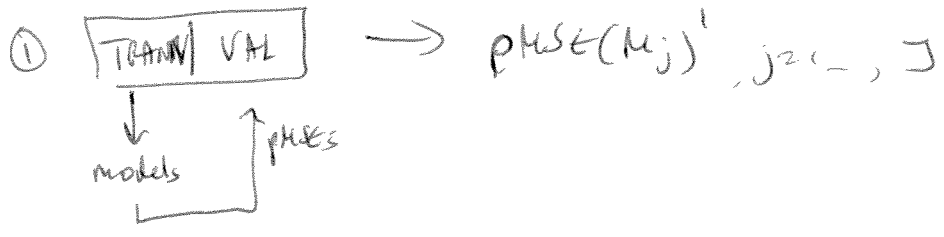
Focus on the

TRAIN	VALIDATION
-------	------------

 part

--	--

 is called 2-fold CV

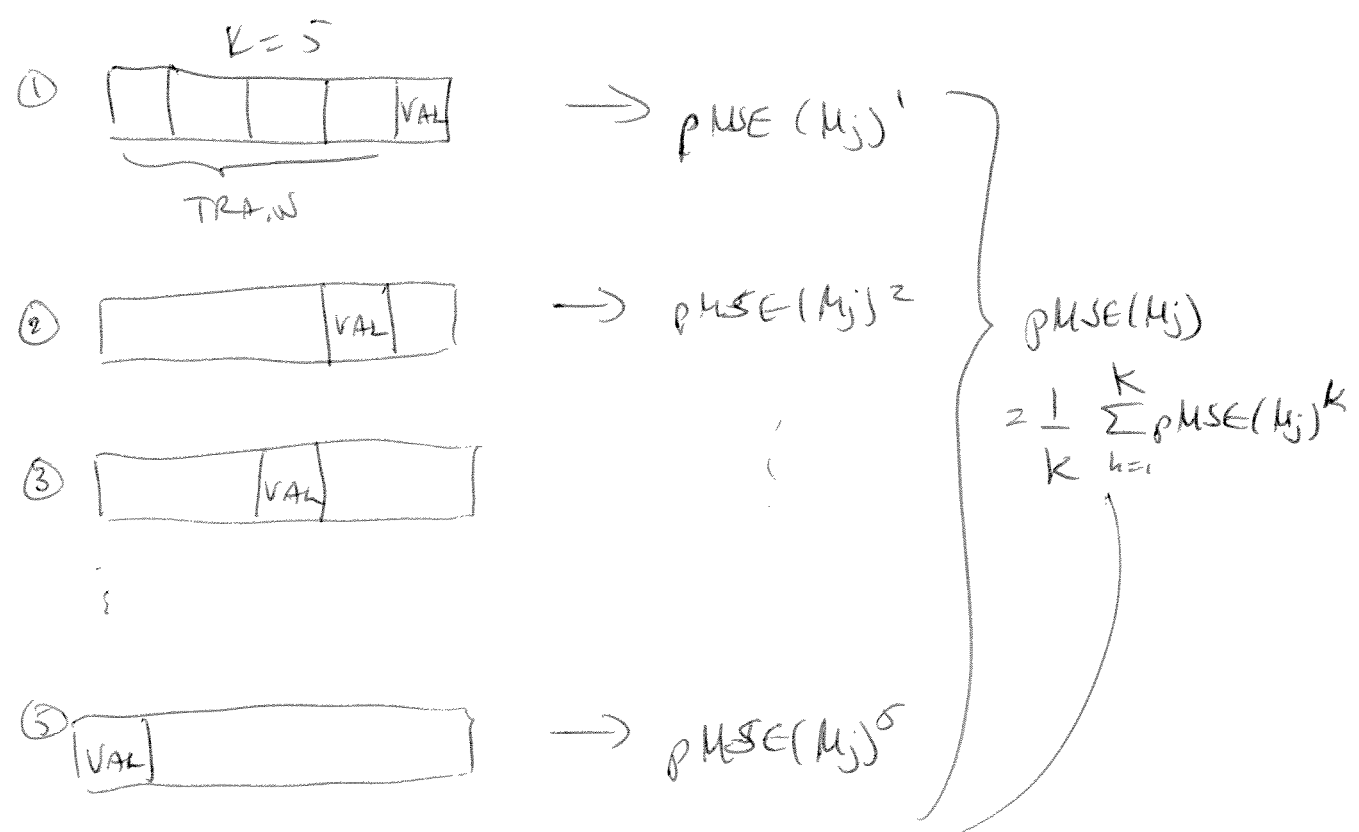


That is, both parts of the data take a turn at training & validation.

We base our final ranking of models on the average

$$PMSE(M_j) = \frac{PMSE(M_j)^1 + PMSE(M_j)^2}{2}$$

Can also do K-fold



Average $PMSE$ on each of the K VALIDATION sets.

How large should K be?

(4)

Small K	Large K
<ul style="list-style-type: none">• TRAINING data much smaller than original data→ more difficult to get good parameter estimates→ bias against complex models <p>(but)</p> <ul style="list-style-type: none">• computationally fast	<ul style="list-style-type: none">• TRAIN almost same size as original data, so almost no bias <p>(but)</p> <p>since each TRAIN share a lot of observations with other TRAIN \Rightarrow individual $PMSE(k_j)^k$ estimates are correlated \Rightarrow high variance estimate of $PMSE$'s</p> <p>→ ranking highly variable too</p>

So small $K \rightarrow$ may bias against selecting big models

large $K \Rightarrow$ $PMSE$'s are highly variable so selection is highly variable too.

Common K 's used — $K=3$
 $K=10$
and $K=n$ (leave-one-out CV)

LOOCV (leave-one-out)

⑤

• Let each observation take turn to be the VARIATION set.

$$\bullet \text{ PMSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}^{-i}(x_i))^2$$

$\underbrace{\hspace{10em}}$
(fitted value @ x_i when (x_i, y_i)
not used for estimation)

Studentized residual $(y_i - \hat{f}^{-i}(x_i))$ can be obtained without refitting the model

$$\hat{f}(x_i) = \sum_{j=1}^n h_{ij} y_j \quad (\text{linear model})$$

$$\Rightarrow \hat{f}^{-i}(x_i) = \sum_{j \neq i} h_{ij} y_j \quad \Rightarrow \hat{f}^{-i}(x_i) = \sum_{j \neq i} h_{ij} y_j + h_{ii} \hat{f}^{-i}(x_i)$$

$\underbrace{\hspace{10em}}$
multiply up $(1 - h_{ii})$

$$\begin{aligned} \Rightarrow (y_i - \hat{f}^{-i}(x_i))^2 &= (y_i - \sum_{j \neq i} h_{ij} y_j - h_{ii} \hat{f}^{-i}(x_i))^2 \\ &= (y_i - \underbrace{\sum_{j \neq i} h_{ij} y_j}_{\hat{f}(x_i)} + h_{ii} (y_i - \hat{f}^{-i}(x_i)))^2 \end{aligned}$$

$$\Rightarrow y_i - \hat{f}^{-i}(x_i) = y_i - \hat{f}(x_i) + h_{ii} (y_i - \hat{f}^{-i}(x_i))$$

$$\Rightarrow y_i - \hat{f}^{-i}(x_i) = \frac{y_i - \hat{f}(x_i)}{1 - h_{ii}}$$

6

$$\Rightarrow \text{LOOCV} = \frac{1}{n} \sum (y_i - \hat{f}^{-i}(x_i))^2$$

$$= \frac{1}{n} \sum \left(\frac{y_i - \hat{f}(x_i)}{1 - h_{ii}} \right)^2$$

So no refitting needed to get the leave-one-out CV.

BUT, for some types of model, hard to compute all leverage values h_{ii}

\Rightarrow Approximation

$$\text{GCV} = \frac{1}{n} \sum \left(\frac{y_i - \hat{f}(x_i)}{1 - \frac{\text{tr}(\hat{H})}{n}} \right)^2$$

replace h_{ii} by average $\frac{\text{tr}(\hat{H})}{n} = \frac{\sum_{j=1}^n h_{jj}}{n}$

For many models $\text{tr}(\hat{H})$ is very easy to compute.

E.g. linear regression $\Rightarrow \text{tr}(\hat{H}) = p$, # variables

2nd approx

Note $\frac{1}{(1-z)^2} \approx 1 + 2z$ if z is small. Here $z = \frac{\text{tr}(\hat{H})}{n}$ small if n is large

$$\Rightarrow \text{GCV} = \frac{\frac{1}{n} \sum (y_i - \hat{f}(x_i))^2}{\left(1 - \frac{\text{tr}(\hat{H})}{n}\right)^2} \approx \underbrace{\frac{1}{n} \sum (y_i - \hat{f}(x_i))^2}_{\text{MSE}} + 2 \cdot \frac{\text{tr}(\hat{H})}{n} \underbrace{\frac{1}{n} \sum (y_i - \hat{f}(x_i))^2}_{\hat{\sigma}^2}$$

$$= \text{MSE} + 2 \cdot p \cdot \hat{\sigma}^2 = C_p!$$

So C_p an estimate of the LOOCV almost