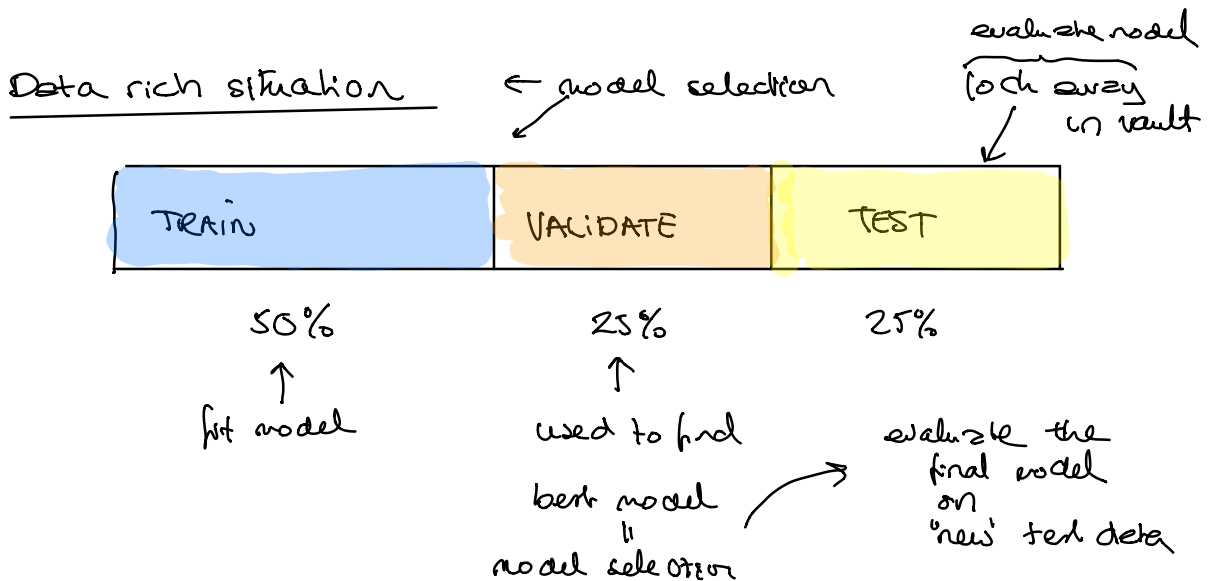


Model selection: choose between different models or different complexity of one model (e.g. KNN).

Model assessment: evaluate performance of chosen model on new data.

Criteria: regression  $E[(y - \hat{f}(x))^2]$  ← MSE <sup>quadratic loss</sup>  
 classification 0/1 loss ← misclassification rate

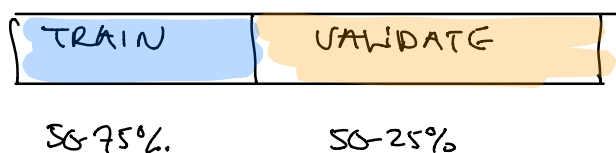


Can we do this in a more data-efficient way?

## CROSSVALIDATION

→ keep the test set in the vault → focus on model selection with the rest of the data.

1) Validation set approach



Ex:  $Y = \text{miles per gallon}$

$X = \text{horsepower} \leftarrow \text{poly } 1, \dots, 10$

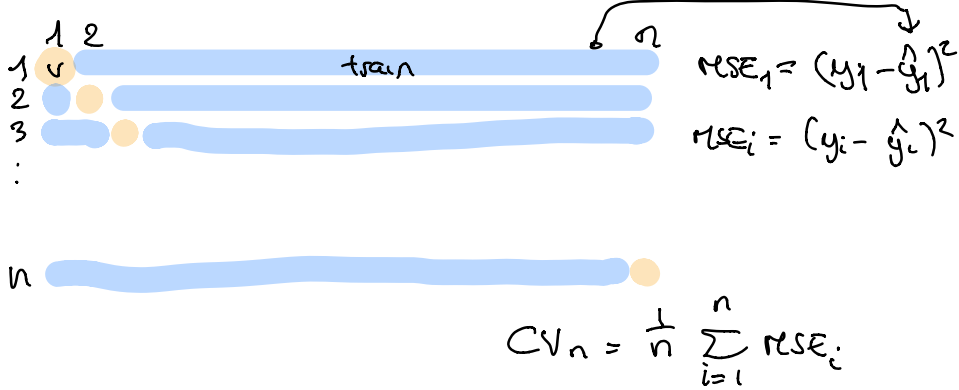
} fit train data for 10 models, use validation set to calc MSE

smallest validation MSE for poly 7

→ the result is dependent on the split of the data into a train/validation set.

→ if we have few data ⇒ exists better solutions!

2) LOOCV leave-one-out-cross-validation



†: no randomness in splits, sample size (n-1)

÷ expensive + high variance

Is  $\frac{n}{k}$  an integer?  $n_1 = n_2 = \dots = n_k$

3) k-fold CV

Shuffle data, make k folds. Draw k=5.

	1	2	3	4	5
1	$n_1$	...	...	...	...
2	...	$n_2$	...	...	...
3	...	...	...	...	...
4	...	...	...	...	...
5	...	...	...	...	$n_5$

from train data

$$MSE_1 = \frac{1}{n_1} \sum_{i \in C_1} (y_i - \hat{y}_i)^2$$

$MSE_2$

$C_1 =$  indices for obs. in fold 1

different fold

want  $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

$$CV_5 = \frac{1}{n} \sum_{j=1}^5 n_j MSE_j$$

"often do with  $\frac{1}{5} \sum_{j=1}^5 MSE_j$ "

after  $n_1 = n_2 = \dots = n_k$  the folds are of nearly eq size 3

$k=5$  and  $k=10$  are very popular, and recommended solutions.

How to choose the best model?

$\theta$  = some parameter to give the model

\* Smallest CV error:

$\theta$  = model parameter

poly:  $\theta=7$  is best for 5 fold

$$\hat{\theta} = \operatorname{argmin} CV(\theta)$$

\* One standard error rule

$$n_1 = n_2 = \dots = n_5$$

MSE<sub>1</sub>, ..., MSE<sub>5</sub>:  $CV_5 = \frac{1}{n} \sum_{j=1}^5 n_j MSE_j$

$$\approx \frac{1}{5} \sum_{j=1}^5 MSE_j$$

$\uparrow$   
var(MSE)

$$SE(\hat{\theta}) \approx (SE(CV_5)) \approx \left(\frac{1}{5}\right)^2 \sum_{j=1}^5 \left(\frac{1}{4} \sum_{j=1}^5 (MSE_j - CV_5)^2\right)$$

Choose the simplest model where  $\frac{1}{5}$

estimator for var(MSE<sub>j</sub>'s)

$$CV(\theta) \leq CV(\hat{\theta}) + SE(\hat{\theta})$$

$\uparrow$   
min rule

give poly 2 in my example

calculate as

$$\sqrt{\frac{\operatorname{var}(mse)}{5}}$$

in R.

CV for both model selection and model assessment?

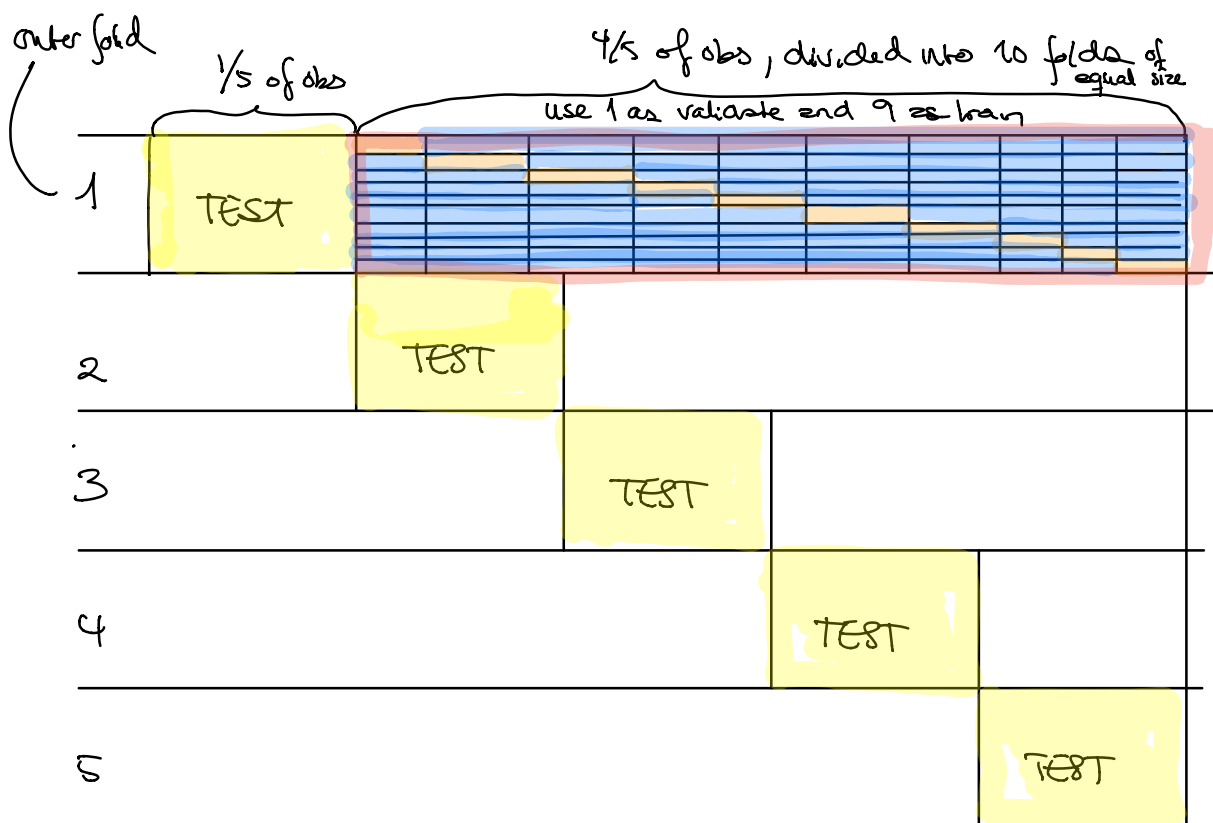
---

→ two nested layers of CV

Outer loop: model assessment  
(draw 5-fold)

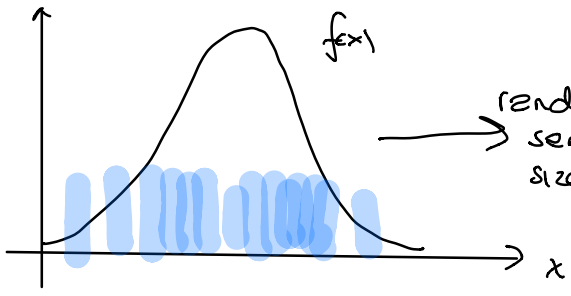
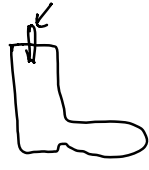
Inner loop: model selection  
(draw 10-fold)

← you choose!



Inner loop: decide on best model → fit to all      and then test on test. ↻ repeat 5 times for the outer loop

# THE BOOTSTRAP



random sample of size  $n$

$x_1, x_2, \dots, x_n$

median

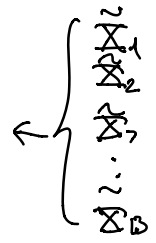
$$\tilde{x}$$



How can I find  $SD(\tilde{x})$

Suggestion: We can draw many <sup>random</sup> samples,  $b=1, \dots, B$  and calculate  $\tilde{x}_b$  for each sample end

$$SD(\tilde{x}) = \sqrt{\frac{1}{B-1} \sum_{b=1}^B \left( \tilde{x}_b - \frac{1}{B} \sum_{b=1}^B \tilde{x}_b \right)^2}$$



[Oh, but in general I have one sample end]  
[I do not know  $f$ ]

New = bootstrapping  $\hat{f}(x) = \frac{1}{n}$  for each  $x_1, x_2, \dots, x_n$

same strategy as before, but now the  $B$  samples are drawn from  $\hat{f}$ .