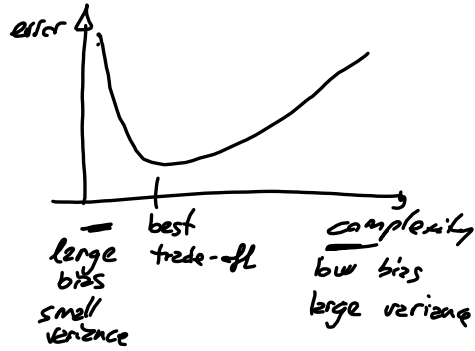


Evaluate the model / find the best complexity

$$E[(f(x) - \hat{f}(x))^2]$$

bias² + variance



- we do not know $f(x)$
- find a way to estimate the expected prediction error with the data we have

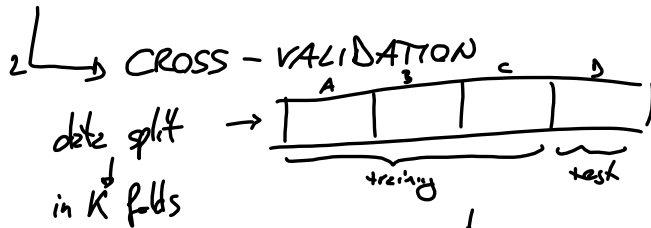
1) random split in training and test set

training	test
50%	50%
66.6%	33.3%
75%	25%

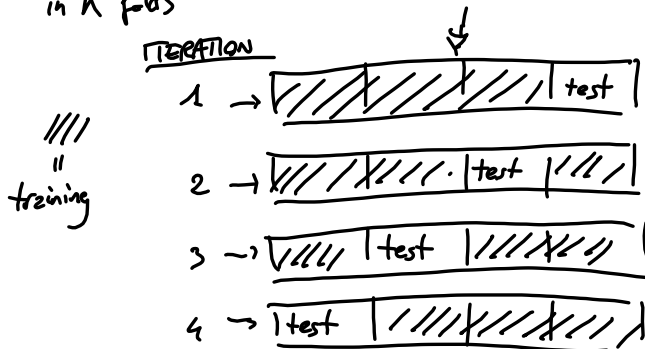


→ when we have 0-1 response, what happens if the response is very unbalanced → stratified ~~resampling~~
 → perform the ~~random~~ split separately for case (1) and controls (0)

- possible issues:
 - it can perform badly if we do not have enough data (not really common in data science, where usually the problem is that the datasets are too big);
 - if we only perform a split, we are not fully using all the information in the data (we do not use the test set to any training part) (we only evaluate on a specific subset (test set) of the data)



$$K=4 \quad \text{err} = \frac{1}{\# \text{ test}} \sum_{i \in \text{test}} (y_i - \hat{f}_{\text{train}}(x_i))^2$$



$$\text{err}_1 = \frac{1}{|D|} \sum_{i \in D} (y_i - \hat{f}^{A,B,C}(x_i))^2$$

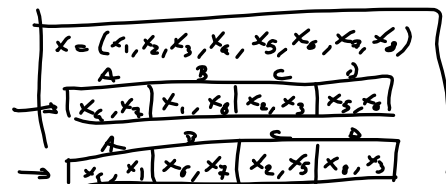
$$\text{err}_2 = \frac{1}{|C|} \sum_{i \in C} (y_i - \hat{f}^{A,B,D}(x_i))^2$$

$$\text{err}_3 = \frac{1}{|B|} \sum_{i \in B} (y_i - \hat{f}^{A,C,D}(x_i))^2$$

$$\text{err}_4 = \frac{1}{|A|} \sum_{i \in A} (y_i - \hat{f}^{B,C,D}(x_i))^2$$

$$\text{err}_{cv} = \frac{1}{4} \sum_{i=1}^4 \text{err}_i$$

y_D	\hat{y}	$\hat{f}^{A,B,C}(x_i)$	
1	0.4	→ 0	1
0	0.7	→ 1	1
1	0.5	→ 1	0
			$\frac{1}{3} (2+1+0) = 0.66$



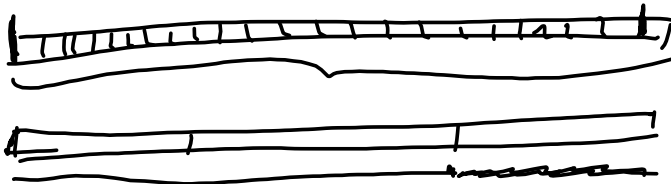
The number of folds is arbitrary (k can be every number $2 \leq k \leq n$), and therefore the procedure is called k -fold cross-validation
 → use in turn 1 of the k folds as a test set, in which we evaluate the models trained on the other $k-1$ folds

Usually we use $k=5$, $k=10$, $k=n$
 ↳ leave-one-out cross-validation

The choice of k also follows a bias-variance trade-off idea

- larger k ⇒ less bias, more variance
- smaller k ⇒ more bias, less variance

→ master course STK-1W4309 we will investigate this in details



→ larger k makes the ^{required} computational time larger (due to increased number of iterations)

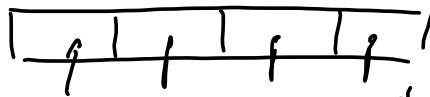
↳ exception: in the case of linear predictor, there is a closed form for the leave-one-out cross-validation

$$\hat{y}_i - y_i = \frac{(y_i - \hat{y}_i)}{1 - p_i} \Rightarrow \text{err}_{\text{loocv}} = \frac{1}{n} \sum_{i=1}^n \frac{(y_i - \hat{y}_i)^2}{(1 - p_i)^2}$$

\uparrow computed $f(x_i)$ \uparrow all points, including i
 \uparrow point set i

Note:

- loocv is a deterministic procedure
- the other k -fold procedure are random (because they are based on a random split)
 → for a practical/computational point of view, set the seed (in R, set.seed(42)) in order to obtain reproducible results.



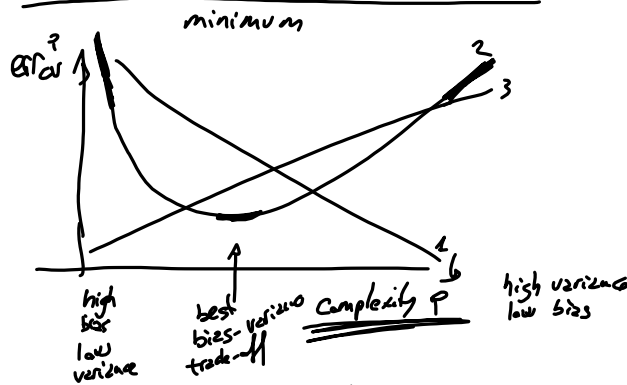
→ repeated cross-validation

- repeat B times
- 1) split observations in k folds
 - 2) perform k -fold cross-validation $\text{err}_{\text{cv}}^{(b)} = \frac{1}{k} \sum_{i=1}^k \text{err}_i^{(b)}$
- average the results $\text{err}_{\text{cv}} = \frac{1}{B} \sum_{b=1}^B \text{err}_{\text{cv}}^{(b)}$

When we want to determine the "best" complexity for a model
 → split data in k -fold
 → train all models ($\text{poly}(x)$, $\text{poly}(x^2)$, ...) on $k-1$ folds
 → test all of them in the remaining fold

test set	$p=1$	$p=2$	$p=3$...	$p=20$
fold 1	$\text{err}_1^{p=1}$	$\text{err}_1^{p=2}$			
fold 2	$\text{err}_2^{p=1}$	\vdots			
\vdots					
fold k	$\text{err}_k^{p=1}$				
average	$\text{err}_{av}^{p=1}$	$\text{err}_{av}^{p=2}$...		$\text{err}_{av}^{p=20}$

$\leftarrow p = \underset{p}{\text{argmin}} \text{err}_{av}^p$



Alternative to cross-validation: bootstrapping

based on non-parametric bootstrap:

1) generate a bootstrap sample, by sampling with replacement n observations from the original sample

e.g. $x = \{x_1, x_2, x_3, x_4, x_5\} \rightarrow x^* = \{\underline{x}_2, x_3, \underline{x}_3, x_4, \underline{x}_2\}$

2) fit the model on the bootstrap sample $\rightarrow \hat{f}^*(x)$

3) use the observation not included on the bootstrap sample as a test set
 e.g. $x_{\text{test}} = \{x_1, x_5\}$

e.g. $\text{err}_b = \frac{1}{2} [(\hat{y}_1 - \hat{f}^*(x_1))^2 + (\hat{y}_5 - \hat{f}^*(x_5))^2]$

4) repeat 1 to 3 B times and average the results

$\text{err}_{\text{boot}} = \frac{1}{B} \sum_{b=1}^B \text{err}_b$

advantage: our training set has the same sample size of the dataset
 → theoretical properties that we will see in STAT 4300

note: bootstrapping is overestimating the error

→ more complex procedures (variants) that take into account this

- 0.632 bootstrapping
- 0.632+

It is important to use as a test set only observations that do not belong to the bootstrap sample

1) generate bootstrap sample x^*

2) fit the model $\hat{f}^*(x)$

3) compute the error on the dataset $\text{err}_b = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}^*(x_i))^2$

4) repeat B times and average errors

- 1) split in training and test
- 2) cross-validation / bootstrapping
- 3) criteriz based on information

We saw:
 - the typical statistical approach for parameter estimation is based on the maximization of the log-likelihood
 - we cannot directly rely on log-likelihood maximization to select the right model \rightarrow more complex model, higher log-likelihood \rightarrow overfitting

Idea:
 - add a penalty term to penalize for too complex model
 \rightarrow not maximize $\underset{\text{(minimize)}}{2 \text{ log-likelihood}}$, but a penalized log-likelihood
 $\rightarrow IC = -2 \text{ log-likelihood} + \text{penalty}(p)$
 penalty(p) penalizes for too complex models, and depends obviously on the complexity

We want to minimize IC.

How to specify penalty(p)

- \rightarrow must be positive and increase with p;
- \rightarrow the penalty for p parameters must be larger than p

$$2(\ell(\hat{\theta}_{p+1}) - \ell(\hat{\theta}_p)) \stackrel{H_0}{\sim} \chi^2_1 \quad H_0: \theta_{p+1} = 0$$

θ_{p+1} is a parameter of dimension p+1
 θ_p " " " " " p

$$E[\chi^2_1] = 1$$

By adding a term without any significance to the model with θ_p , we expect a decrease in terms of $-2\ell(\theta)$ of 1

$$IC = -2 \ell(\hat{\theta}) + \text{penalty}(p)$$

Typical choices for penalty(p) are

NAME	penalty(p)
Akaike Information Criterion (AIC)	$2p$
corrected AIC (AIC _c)	$2p + \frac{2p(p+1)}{n-(p+1)}$
Bayesian Information Criterion (BIC)	$p \log(n)$

AIC

Idea: we want to minimize the Kullback-Leiber divergence between the true model and the fitted model

$$KL(p_0(y), p(y, \theta)) = E_p \left[\log \frac{p_0(y)}{p(y, \theta)} \right] = \underbrace{E_p [\log p_0(y)]}_{\text{true model}} - \underbrace{E_p [\log p(y, \theta)]}_{\text{fitted model}}$$

independent of θ \rightarrow we can only act on this part \Rightarrow the largest, the better
 \Downarrow
 maximum likelihood

The key point is that the estimation of $E_p [\log p(y, \theta)]$ is

$$\log p(y, \theta) - p$$

multiplied by the usual factor -2 $\rightarrow AIC = -2 \log(p(y, \theta)) + 2p$
 \hookrightarrow align that to the standard likelihood quantities

Alternative interpretation (Hastie, Tibshirani & Friedman, 2009) : $-2 \log p(y; \hat{\theta}) + 2p$ correction for (over)optimism
 estimation of the training error STK-W4300 why
 $AIC = \text{training error} + \hat{w}$

AIC_c is a correction for small sample size

$$2p + \frac{2p(p+1)}{n-(p+1)} \xrightarrow{n \rightarrow \infty} 2p$$

When $n \rightarrow \infty$, $AIC = AIC_c$

BIC is motivated differently (STK-W4300)

Consider the two penalties

AIC
 $2p$

BIC
 $p \log(n)$

$\log n > 2 \quad n > e^2 \approx 7.4$

less complex model for $n \geq 8$

AIC: better for prediction (when including a non-relevant variable is not too bad) as long as the variance is not increased too much

BIC: support less complex (sparser) models.

In theory, BIC works the best (is consistent)

In practice, AIC works better (BIC tends to select too sparse models)