

MA8701 Advanced methods in statistical inference and learning

~~L3-4~~: Model selection and assessment

L4:

Mette Langaas

~~1/17/23~~ 2023.01.20

Cross-validation (CV)

(ESL Ch 7.10, 7.12 - most should be known from TMA4268)

The aim is to estimate Err_T , but from simulation analyses (ESL Ch 7.12) it turns out that cross-validation estimates Err “the best”.

The starting point for the method is that we only have one training set - and try to use that for either model selection or model assessment (not both).

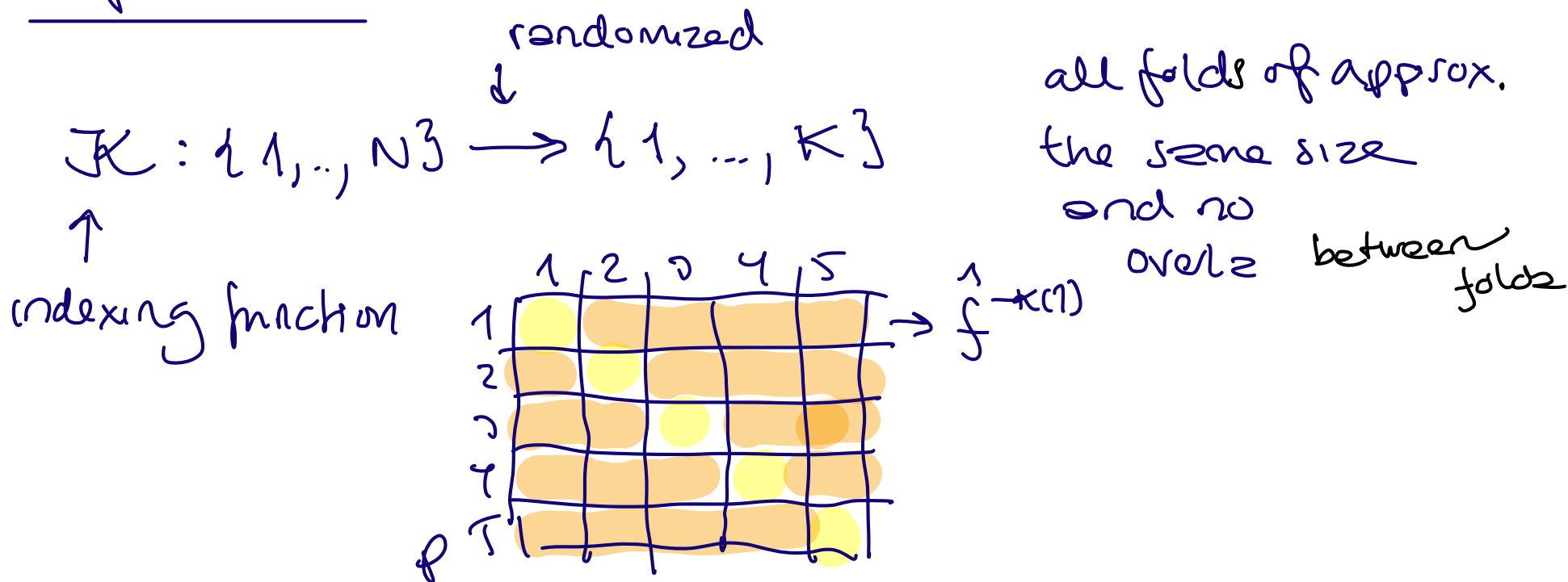
What to do when both is needed, is not covered in this chapter. Nested cross-validations aka two-layers of cross-validation is one possibility. Another is to set aside data for a test set for model assessment, but use the training set in cross-validation for model selection.

CROSS-VALIDATION (ESL 7.10)

We will estimate $\text{Err} = E_{X,Y} (L(Y, \hat{f}(X)))$ extra sample error

NB: Err, not Err_n !

K-fold CV



$\hat{f}_{-k}(x)$: estimator based
 on all data except on fold k

Err:

$$CV(\hat{f}^\alpha) = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{f}_{-k(i)}^\alpha(x_i))$$

tuning parameter α

the fold where x_i was
 not a part of the
 "training set"

Model selection: choose α where $CV(\hat{f}^\alpha)$ is minimum or
 use 1-SD-rule

Formal set-up for model assessment

- ▶ The allocation of observation $\{1, \dots, N\}$ to folds $\{1, \dots, K\}$ is done using an indexing function $\kappa : \{1, \dots, N\} \rightarrow \{1, \dots, K\}$, that for each observation allocate the observation to one of K folds.
- ▶ Further, $\hat{f}^{-k}(x)$ is the fitted function, computed on the observations except the k th fold (the observations from the k th fold is removed).
- ▶ The CV estimate of the expected prediction error $\text{Err} = \mathbb{E}_T \mathbb{E}_{X^0, Y^0} [L(Y^0, \hat{f}(X^0)) \mid T]$ is then

$$\text{CV}(\hat{f}) = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}^{-k(i)}(x_i))$$

Efron and Hastie (2016) page 218: “ $\text{CV}(\hat{f})$ is estimating the average prediction error of the algorithm producing \hat{f} , not \hat{f} itself”.

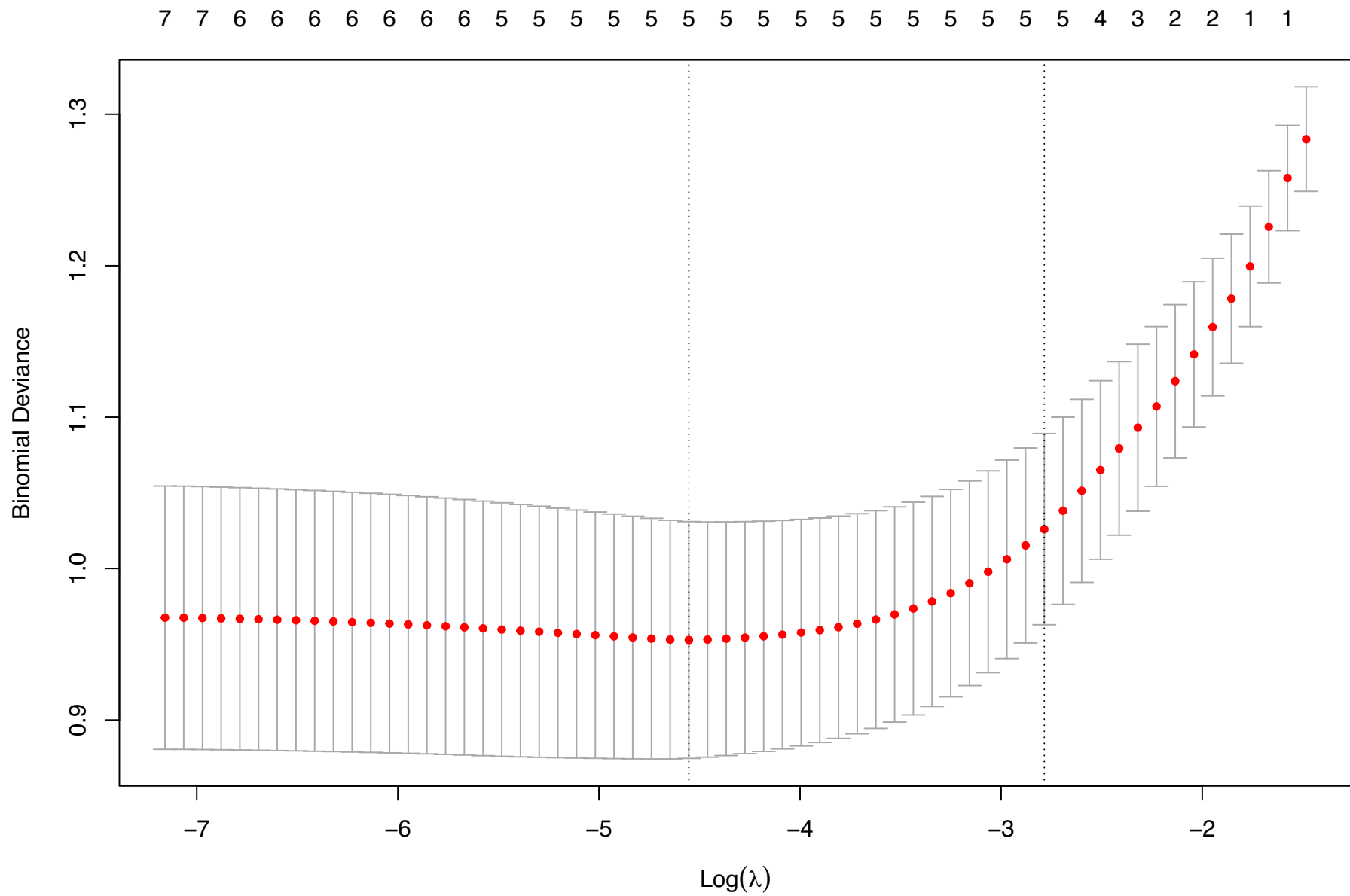
Pima indian example

We will use the classical data set of *diabetes* from a population of women of Pima Indian heritage in the US, available in the R MASS package. This version of the data has no missing values. The following information is available for each woman:

- ▶ diabetes: 0= not present, 1= present
- ▶ npreg: number of pregnancies
- ▶ glu: plasma glucose concentration in an oral glucose tolerance test
- ▶ bp: diastolic blood pressure (mmHg)
- ▶ skin: triceps skin fold thickness (mm)
- ▶ bmi: body mass index (weight in kg/(height in m)²)
- ▶ ped: diabetes pedigree function.
- ▶ age: age in years

We will use the default division into training and test in the MASS library, with 200 observations for training and 332 for testing.

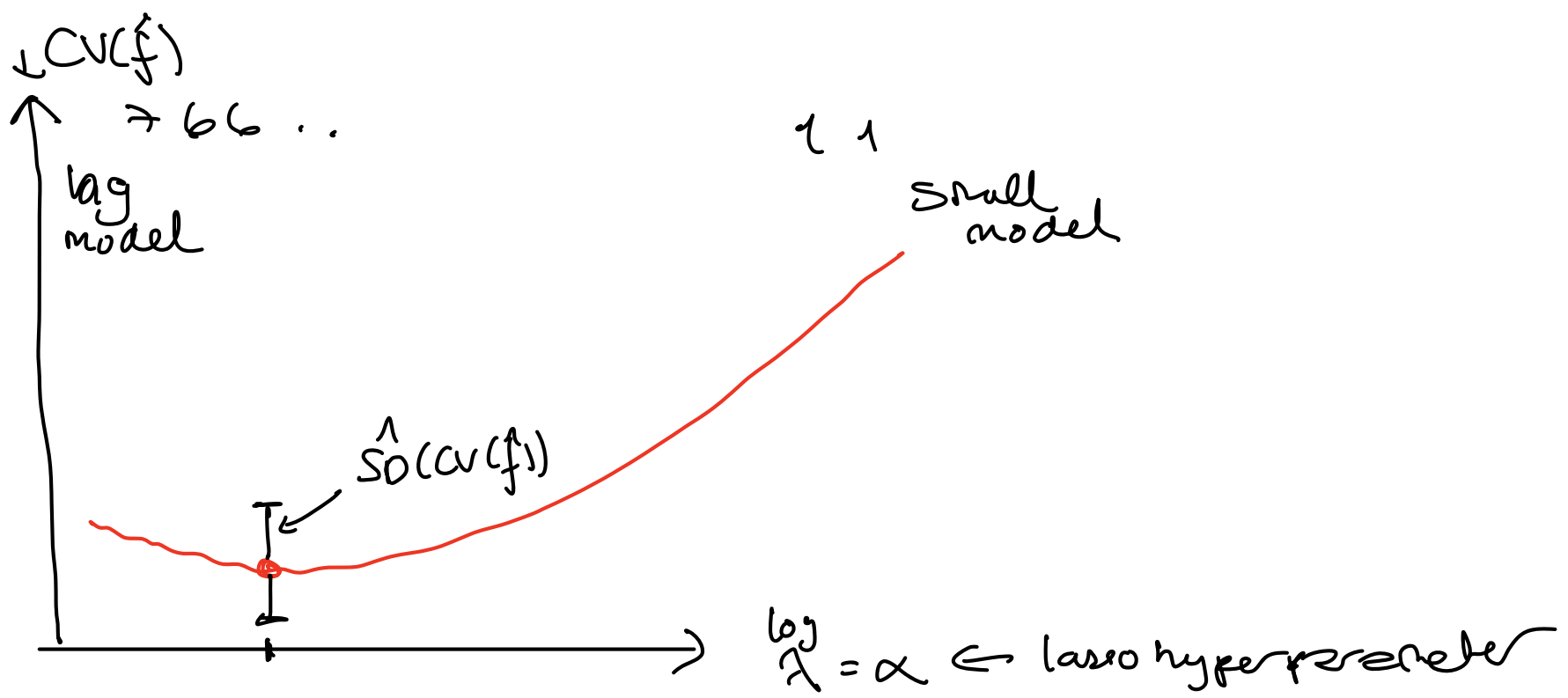
Only training set used here.



Group discussion

The lasso logistic regression (to be studied in Part 2) was used to fit the data, and some loss function is plotted on the vertical axis (more in Part 2) and on the horizontal axis the loss for different fits for different choices of a complexity parameter is given. 10-fold crossvalidation is used. (Just assume that a generic prediction is used, this is not meant to be specific for the lasso.)

- ▶ What are the red dots and how have they been calculated?
- ▶ What are the the vertical bars sticking out of each red dot, and how have they been calculated? What do they picture?
- ▶ What (your choice) is the optimal choice of the complexity parameter?



Red dot: $CV(\hat{f}(x_i))$, Vertical bar: $\hat{S}_D(CV(\hat{f}))$

$$\hat{Err} = CV(\hat{f}) = \frac{1}{N} \sum_{i=1}^N h(y_i, \hat{f}^{-k(i)}(x_i)) = \frac{1}{K} \sum_{j=1}^K \hat{CV}_j$$

$$\frac{1}{N_j} \sum_{i \in K_j(i)} h(y_i, \hat{f}^{-k(i)}(x_i))$$

$$\text{Want: } \text{Var}(\hat{\text{Err}}) = \text{Var}\left(\frac{1}{n} \sum_{j=1}^k \hat{C}_j^a\right) = \frac{1}{n^2} \sum_{j=1}^k \text{Var}(\hat{C}_j^a)$$

$$= \frac{1}{n} \underbrace{\text{Var}(\hat{C}_j^a)}$$

$$\hat{\text{Var}}(\text{CV}(\hat{f})) = \frac{1}{n} \cdot \frac{1}{k-1} \sum_{j=1}^k (\hat{C}_j^a - \text{CV}(\hat{f}))^2$$

$\hat{\text{SD}}(\text{CV}(\hat{f}))$ is shown as vertical bar

"we choose the most parsimonious model whose error is no more than one sd above the error of the best model"

1 SD-error rule

Choice of K

- ▶ Popular choices are 5 and 10 based on observations in simulation studies- and arguments similar to a bias-variance trade off.
- ▶ $K = N$ is called *leave-one-out* cross-validation LOOCV, and gives the lowest bias for estimating the Err.

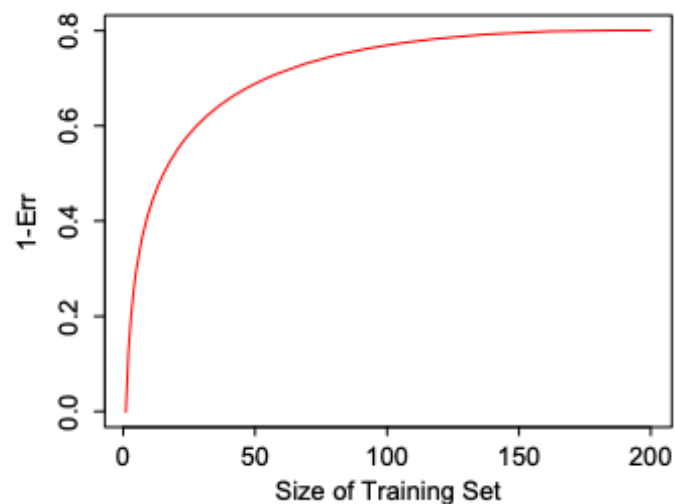


FIGURE 7.8. *Hypothetical learning curve for a classifier on a given task: a plot of $1 - \text{Err}$ versus the size of the training set N . With a dataset of 200 observations, 5-fold cross-validation would use training sets of size 160, which would behave much like the full set. However, with a dataset of 50 observations fivefold cross-validation would use training sets of size 40, and this would result in a considerable overestimate of prediction error.*

Generalized cross-validation (GCV)

For LOOCV with squared loss and linear fitting. Remember

$$\hat{\mathbf{y}} = \mathbf{S}y$$

For many fitting methods (including MLR)

$$\frac{1}{N} \sum_{i=1}^N [y_i - \hat{f}^{-i}(x_i)]^2 = \frac{1}{N} \sum_{i=1}^N \left[\frac{y_i - \hat{f}(x_i)}{1 - S_{ii}} \right]^2$$

where S_{ii} is the i th diagonal element of \mathbf{S} . This leads to the GCV approximation:

$$\text{GCV}(\hat{f}) = \frac{1}{N} \sum_{i=1}^N \left[\frac{y_i - \hat{f}(x_i)}{1 - \text{tr}(\mathbf{S})/N} \right]^2$$

where we recognise the effective number of parameters $\text{trace}(\mathbf{S})$. In some settings the $\text{trace}(\mathbf{S})$ is computed more easily than the individual elements S_{ii} .

The wrong and the right way to do cross-validation

In short: make sure that all part of the model fit process is “inside” the CV.

See learning material from TMA4268: Module 5: Resampling, and I also recommend to work on Problem 3 with solutions

Group discussion

Can you give one example of a right way to do cross-validation and also a wrong way? If you want you may use the Pima-indians as an example, but other examples may also be used.

- randomized allocation to folds, maybe also stratified
- remember also standardization part of CV-loop

Selection bias in gene extraction on the basis of microarray gene-expression data

Article by [Christophe Ambroise and Geoffrey J. McLachlan, PNAS 2002](#): Direct quotation from the abstract of the article follows.

- In the context of **cancer diagnosis and treatment**, we consider the problem of constructing an accurate prediction rule on the basis of a relatively small number of tumor tissue samples of known type containing the expression data on very many (possibly thousands) genes.
- Recently, results have been presented in the literature suggesting that it is **possible to construct a prediction rule from only a few genes such that it has a negligible prediction error rate**.
- However, in these results the test error or the leave-one-out cross-validated error is **calculated without allowance for the selection bias**.

-
- There is no allowance because the rule is either tested on tissue samples that were used in the first instance to select the genes being used in the rule or because the cross-validation of the rule is not external to the selection process; that is, gene selection is not performed in training the rule at each stage of the crossvalidation process.
 - We describe how in practice the selection bias can be assessed and corrected for by either performing a crossvalidation or applying the bootstrap external to the selection process.
 - **We recommend using 10-fold rather than leave-one-out cross-validation, and concerning the bootstrap, we suggest using the so-called .632 bootstrap error estimate designed to handle overfitted prediction rules.**
 - Using two published data sets, we demonstrate that when correction is made for the selection bias, the cross-validated error is no longer zero for a subset of only a few genes.

0.632+ ?

BOOTSTRAP METHODS (ESL 7.11)

We study nonparametric bootstrap - sampling from data - and focus on estimating Err .

$$Z_1 = (x_1, y_1), Z_2 = (x_2, y_2), \dots, Z_N = (x_N, y_N)$$

are independent observation pairs from the training set

$$Z = \{Z_1, \dots, Z_N\}$$

↓ draw with replacement N obs. pairs Z^{*b} ← b th sample

↓ repeat B times

$$Z^{*1}, Z^{*2}, \dots, Z^{*B}$$

↓
calculate the quantity of interest: \hat{f}^{*b} for each boot sample

How to use \hat{f}^{*b} to estimate Err ?

First try: $\hat{\text{Err}}_{\text{boot}} = \frac{1}{B} \frac{1}{N} \sum_{b=1}^B \sum_{i=1}^N L(y_i, f^{*b}(x_i))$

Z^{*b} : act like a training set

$Z_{\text{arg data}}$: act like test set

Quality of $\hat{\text{Err}}_{\text{boot}}$? Too low \nwarrow $\hat{\text{Err}}_{\text{boot}}$ too low!

Because overlap between
"train" and "test" data!

However - $\widehat{\text{Err}}_{\text{boot}}$ is not a good estimator: bootstrap datasets are acting as training data and the original data as a test sample - and the two samples have observations in common.

This overlap can make predictions too good. Remember, in CV we have no overlap.

Q: What is the probability that observation i is included in bootstrap sample b ?

In one draw:

$$P(\text{draw } x_i) = \frac{1}{N}, \quad P(\text{not draw } x_i) = 1 - \frac{1}{N}$$

In N_{\uparrow} draws: $P(\text{not draw } x_i \text{ in any draw}) = \left(1 - \frac{1}{N}\right)^{N_{\uparrow}}$
independent

$$P(x_i \text{ is included in boot sample } b) = 1 - (1 - \frac{1}{n})^n$$

$$\xrightarrow{n \rightarrow \infty} \approx 1 - e^{-1} = 0.632$$

$$P(\text{not included}) = 0.368$$

Ex: two-class problem where x and y are independent
Err = 0.5 if use

0/1
↓
↑ random assignment

↑
1-NN method: obs i is in boot b : $h(y_i)_{b}(x_i) = 0$
else random guess

$$0.632 \cdot 0 + 0.368 \cdot 0.5 = 0.184 \quad \text{too low!}$$

Second try: leave-one-out bootstrap aka
out of bag OOB

$$\hat{\text{Err}}^{(1)} = \frac{1}{N} \sum_{i=1}^N \frac{1}{|C^{-i}|} \sum_{b \in C^{-i}} L(y_i, \hat{f}^{ab}(x_i))$$

how many such
 boot samples where z_i not included
 no overlap between y_i and f by construction!

\hat{f}^{ob} is based on bootstrap sample b , with N observations, but what is the average number of distinct obs in \hat{z}^b ? (distinct: we assume \hat{z} has N distinct obs)

$$\begin{aligned}
 & E(\text{number of distinct obs drawn}) \\
 &= E(I_1 + I_2 + \dots + I_N) = N \cdot (1 - (1 - \frac{1}{N})^N) \\
 &\approx N \cdot 0.632
 \end{aligned}$$

$$I_i = \begin{cases} 0 & \text{not drawn} \\ 1 & \text{drawn} \end{cases}$$

$$E(I_i) = P(I_i = 1) = 1 - (1 - \frac{1}{N})^N$$

The estimator \hat{f}^{nb} will then be using very little of the ^{original} training data $\Rightarrow \hat{Err}^{(1)}$ too big!

Now - let us use $\hat{Err}^{(1)}$ (too big) together with

$$\overline{err} = \frac{1}{n} \sum_{i=1}^n h(y_i, \hat{f}(x_i)) \text{ (too low)} - \text{in a weighted fashion}$$

Third try:

↑
what is this weight?

$$\hat{Err}^{(1)} (0.632) = 0.368 \overline{err} + 0.632 \cdot \hat{Err}^{(1)}$$

good performance in light fitting situations

but not for "overfitted situations".

Ex: $\overline{err} = 0$ because $h=0$ always in both.

$$\hat{Err}^{(1)} = 0.5 \quad \hat{Err}^{(1)} (0.632) = 0.632 \cdot 0.5 + 0 = 0.316 \quad \text{true } 0.5$$

According to ESL (page 251): the derivation of the .632 estimator is complex, and the estimator is expected to work well in situation where the data is not overfitted, but may break down in overfit situations.

According to CASI (page 323) the .632 rule is less variable than the leave-one-out CV.

Example of this on page 251-252: two equal size classes where predictors independent of class, classification with 1NN gives $\overline{\text{err}} = 0$ and $\widehat{\text{Err}}^{(1)} = 0.5$ and thus $\widehat{\text{Err}}^{(.632)} = 0.632 \cdot 0.5 = 0.316$, where here the true error rate is 0.5.

Fourth try:

$$\hat{\text{Err}}^{(0.632+)} = (1 - \hat{w}) \overline{\text{err}} + \hat{w} \hat{\text{Err}}^{(1)}$$

\hat{w} is decided by evaluating the amount of
(overfitting) overfitting.

There is an improved version of the estimator - taking into account the amount of overfitting, leading to an adjustment to the weight $w = 0.632$ (and $1 - w = 0.368$) dependent on a so-called *no-information error rate* $= \gamma$ = the error rate of the prediction rule when predictors and class labels are independent.

$$\hat{\gamma} = \frac{1}{N^2} \sum_{i=1}^N \sum_{i'=1}^N L(y_i, \hat{f}(x_{i'}))$$

\swarrow obs i \swarrow prediction at a different $x_{i'}$

Further the *relative overfitting rate* is defined to be

$$\hat{R} = \frac{\widehat{\text{Err}}^{(1)} - \overline{\text{err}}}{\hat{\gamma} - \overline{\text{err}}}$$

$\hat{R} = 0 \Leftrightarrow \widehat{\text{Err}}^{(1)} = \overline{\text{err}}$ so no overfitting and $\hat{w} = 0.632$ used

$\hat{R} = 1 \Leftrightarrow \widehat{\text{Err}}^{(1)} = \frac{1}{2}$ no-info case and $\hat{w} = 1$

Finally, the “.632+”-estimator is

$$\widehat{\text{Err}}^{(.632+)} = (1 - \hat{w})\overline{\text{err}} + \hat{w}\widehat{\text{Err}}^{(1)}$$

where $\hat{w} = \frac{0.632}{1 - 0.368\hat{R}}$.

For details on this approach consult ESL page 252-253.

The 0.632+ bootstrap estimator is a method for estimating the performance of a machine learning model on unseen data, based on its performance on a sample of data. It is an extension of the 0.632 bootstrap estimator, which was proposed by Efron in 1983.

The 0.632 bootstrap estimator is based on the idea that the performance of a model on the training set is likely to be better than its performance on unseen data. To estimate the performance on unseen data, the 0.632 bootstrap estimator takes the average of the performance on the training set and the performance on a validation set, with a weight of 0.632 given to the performance on the training set.

The 0.632+ bootstrap estimator extends this idea by taking into account the variability of the performance on the validation set. Instead of taking a fixed weight of 0.632 for the performance on the training set, the 0.632+ bootstrap estimator estimates the weight dynamically, by taking into account the variability of the performance on the validation set.

The 0.632+ bootstrap estimator is calculated as follows:

- Take a sample of the dataset with replacement, called bootstrap sample.
- Train the model on the bootstrap sample and test it on the remaining data.
- Repeat this process B times.
- Calculate the average performance of the model across all B runs.
- Let the performance of the model on the bootstrap sample be denoted by T and the performance on the remaining data be denoted by V.
- The 0.632+ bootstrap estimator is given by: $T + (V - T) * (1 - B^{-1/2})$

The 0.632+ bootstrap estimator is a more robust method for estimating the performance of a model on unseen data, as it takes into account the variability of the performance on the validation set. It is particularly useful when the sample size is small or the performance of the model is highly variable.

openai ChatGPT:
Explain the 0.632+ bootstrap estimator

Can we think of the bootstrap sample as a validation set? maybe?

0.368 \overline{err}

anything

overfitting

The OOB data is more of a training set than a test set in this setting? The sum in the OOB bootstrap estimator is not this

0
OOB

$\hat{Err}^{(1)}$

$\hat{Err}^{(0.632+)}$

$= (1 - \hat{w}) \overline{err}$

$+ \hat{w} \hat{Err}^{(1)}$

Is B in the formula?

GRADE: E/F?

Conclusions

Group discussion:

Construct a “mind map” or “overview sheet” or “concept map” (mind map with verbs on arrows between entities) for the “Model assessment and selection” topics, and write down important take home messages!

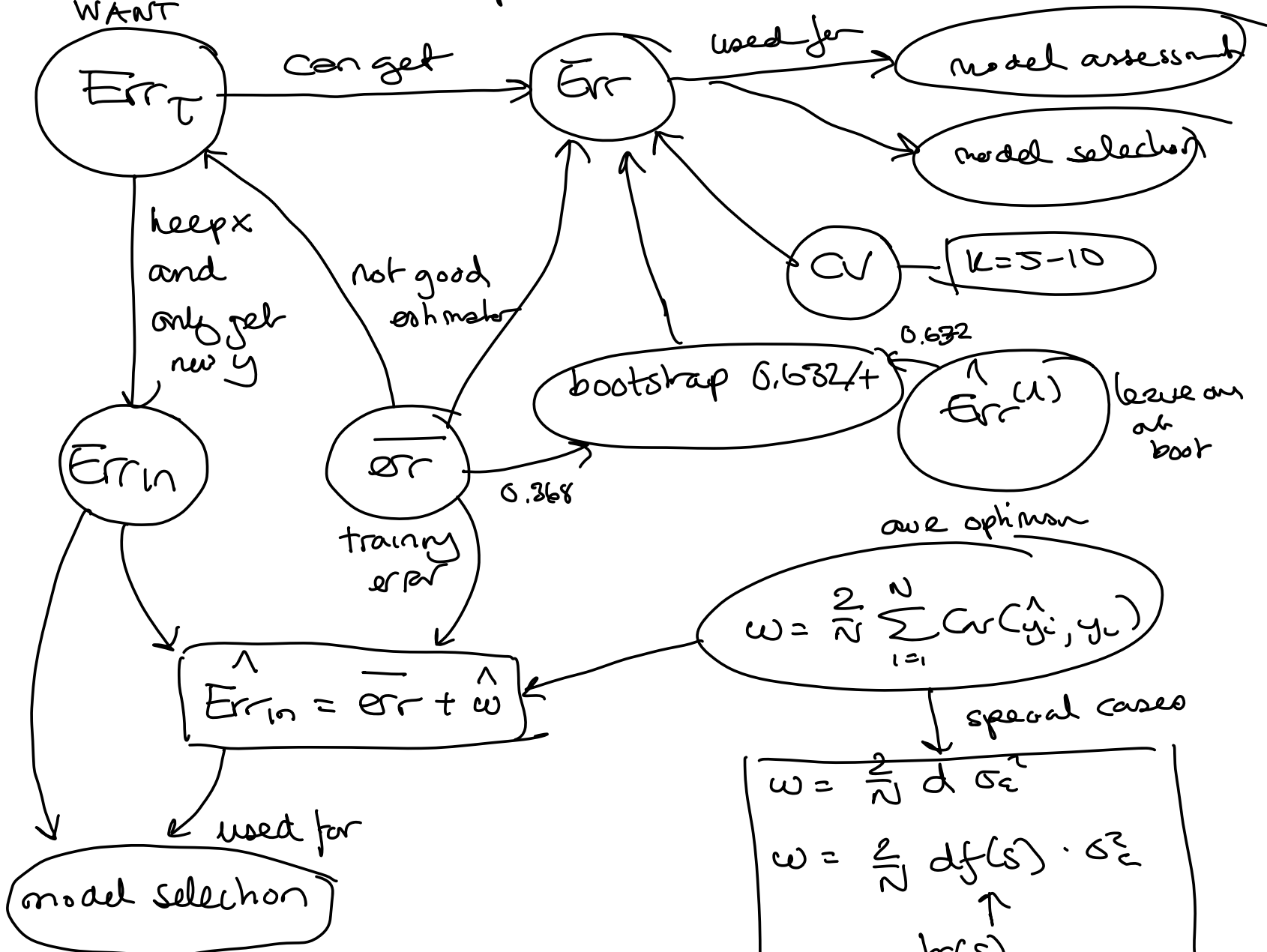
Some concept that could be in the map:

$\text{Err}, \text{Err}_T, \text{Err}_{\text{in}}, \widehat{\text{Err}}_{\text{in}}, \overline{\text{err}}, \omega, \text{Cov}(\hat{y}_i, y_i), \text{trace}(S), \text{df}, \text{CV}(\hat{f}), \text{bootstrap},$
 $\widehat{\text{Err}}^{(1)}, 0.632, 0.368, \text{model assessment}, \text{model selection}.$

Cond. ZPE

EPE = expected prediction error

WANT



$$w = \frac{2}{N} \sum_{i=1}^N Cr(\hat{y}_i, y_i)$$

special cases

$$w = \frac{2}{N} d \sigma_e^2$$

$$w = \frac{2}{N} df(s) \cdot \sigma_e^2$$

↑
h(s)
↑
ŷ = SY

Final remarks

- ▶ In a perfect world we would be rich on data and can divide available data into sets for training, validation and testing
- ▶ We have derived cool covariance-result on expected optimism for training error related to in-sample prediction error (the covariance) - that is used for finding model selection criteria (but not for model assessment). If we can't calculate a formula for the covariance, bootstrapping can be used to do this (Efron and Hastie (2016) Equation 12.64 on page 224).
- ▶ Estimating expected prediction (test) error for a particular training set is not easy in general (if we only have this one training set), but cross-validation and bootstrapping may provide reasonable estimates of the expected test error Err .
- ▶ If resampling needed for model assessment: take average of many 10-fold CV Err estimates?

Week 3: How to handle missing data (new this year)

Reading list: online-book of van Buuren - given chapters and Handbook of missing data (oria) parts of chapter 12 on multiple imputation.

Monday 23.01: Motivation, missing types, different solutions, starting on multiple imputation

Friday 27.01: Mainly multiple imputation

Friday 27.01 after class: meeting with reference group.

Usual agenda for such meetings:

We go through the learning outcome and compare with the learning plan/ resources and evaluation - and see if they are in alignment (samstemt undervisning)

Comment on lectures, exercises, text book, plan for the course in general.

For Part 4: should we be full digital the two times that Kjersti presents from Oslo, or hybrid (then difficult for Kjersti to help in break-out rooms). The third time Kjersti is in Trondheim (3 hrs physical).

Data analysis project 1: based on Part 2. Which deadlines should we have?

Article presentation: During Part 3 and also possibly Part 5 (after Easter).