

Chapter 4

Cross-Validation

4.1 Introduction

The generalization performance (which is a terminology from the machine learning community) of a learning method describes its prediction capacity on new test data. Or in other words, the generalization performance measures the predictive power of a learning method on new, out-sample data. Having an estimate of this predictive power is very important for comparing among different methods or for tuning an estimator or algorithm to optimize predictive power.

Cross-Validation is one very general method for estimating such generalization or out-sample performance.

4.2 Training and Test Set

Consider the situation where we have data

$$(X_1, Y_1), \dots, (X_n, Y_n) \text{ i.i.d. } \sim P$$

where P denotes the unknown distribution. (In case where the x_i is non-random, we can also drop the i.i.d. assumption for the pairs (x_i, Y_i) .)

We typically have a target in mind, for example the regression function $m(\cdot)$. The estimated regression function $\hat{m}(\cdot)$ is then constructed by using some estimator or algorithm, based on the data $(X_1, Y_1), \dots, (X_n, Y_n)$. This data is also called the **training data**, since it is used to train the estimator or learning algorithm.

We then would like to evaluate the accuracy of the estimated target which is based on the training data. A principal problem thereby is that if we use the training data again to measure the predictive power of our estimated target (e.g. of $\hat{m}(\cdot)$), the results will be overly optimistic. For example, when using the residual sum or squares in regression,

$$\sum_{i=1}^n (Y_i - \hat{m}(X_i))^2$$

becomes smaller the more “complex” (more degrees of freedom) the model for $\hat{m}(\cdot)$ involves. Obviously, a very complex model will not be good. We have seen already in Section 3.4 that penalizing the residual sum of squares can be useful to cope with the overfitting problem.

Alternatively, we could look how well the estimated target, e.g. $\hat{m}(\cdot)$, does on **new test data**

$$(X_{new,1}, Y_{new,1}), \dots, (X_{new,\ell}, Y_{new,\ell}) \text{ i.i.d. } \sim P,$$

which is assumed to be independent from the training data but having the same distribution P . In case of regression, we would like to evaluate

$$\ell^{-1} \sum_{i=1}^{\ell} (Y_{new,i} - \hat{m}(X_{new,i}))^2.$$

More generally, when having an estimated target \hat{m} and a loss function ρ , we would like to evaluate

$$\ell^{-1} \sum_{i=1}^{\ell} \rho(Y_{new,i}, \hat{m}(X_{new,i})),$$

where $\hat{m}(\cdot)$ is constructed from training data only. If ℓ is large, this evaluation on the test set approximates the theoretical test set error

$$\mathbb{E}_{(X_{new}, Y_{new})} [\rho(Y_{new}, \underbrace{\hat{m}}_{\text{based on training data only}}(X_{new}))]$$

which is still a function of the training data. The expected value of this (with respect to the training data) is the generalization error,

$$\mathbb{E}_{\text{training}} \mathbb{E}_{(X_{new}, Y_{new})} [\rho(Y_{new}, \underbrace{\hat{m}}_{\text{...training data}}(X_{new}))] = \mathbb{E}[\rho(Y_{new}, \hat{m}(X_{new}))], \quad (4.1)$$

where the latter \mathbb{E} is over the training data $(X_1, Y_1), \dots, (X_n, Y_n)$ as well as the test data (X_{new}, Y_{new}) . This generalization error avoids the fact that it would get smaller the more complex the model for $\hat{m}(\cdot)$. In particular, it will increase as \hat{m} is overfitting. Hence the generalization error in (4.1) is a very useful quantity to regularize and tune an estimator or algorithm for \hat{m} .

4.3 Constructing training-, test-data and cross-validation

Unfortunately, test data is not available at the time we want to run and tune our estimator or algorithm. But we can always “artificially” construct (smaller) training- and test-data.

4.3.1 Leave-one-out cross-validation

For leave-one-out cross-validation (CV), we use the i th sample point as test data (test sample size = 1) and the remaining $n - 1$ sample points as training data.

Denote in general the estimator or algorithm by $\hat{\theta}_n$ which is based on the n sample points. In CV, when deleting the i th sample, the estimator is based on the sample without the i th observation, and we denote this by

$$\hat{\theta}_{n-1}^{(-i)}, \quad i = 1, \dots, n.$$

We can then evaluate this estimate on the i th observation (the test sample), for every $i = 1, \dots, n$. To make this more explicit, we suppose that the estimator $\hat{\theta}$ is a curve

estimator \hat{m} , and performance is evaluated in terms of a loss function ρ , e.g. $\rho(u) = u^2$ as in Chapter 3. The cross-validated performance is then

$$n^{-1} \sum_{i=1}^n \rho \left(Y_i, \hat{m}_{n-1}^{(-i)}(X_i) \right) \quad (4.2)$$

which is an estimate of the test set error, or generalization error, in (4.1).

Note that the CV-method is very general: it can be used for any loss function ρ and in many problems which can be different than linear or nonparametric regression.

CV in general requires that the estimator $\hat{\theta}$ is fitted n -times, namely for all training sets where the i th observation has been deleted ($i = 1, \dots, n$).

4.3.2 K -fold Cross-Validation

A computationally cheaper version of leave-one-out CV is the so-called K -fold CV. The data set is randomly partitioned into K equally sized (as equal as possible) subsets \mathcal{B}_k of $\{1, \dots, n\}$ such that $\cup_{k=1}^K \mathcal{B}_k = \{1, \dots, n\}$ and $\mathcal{B}_j \cap \mathcal{B}_k = \emptyset$ ($j \neq k$). We can now set aside a k th test data set including all sample points whose indices are elements of \mathcal{B}_k .

K -fold cross-validation then uses the sample points with indices not in \mathcal{B}_k as training set to construct an estimator

$$\hat{\theta}_{n-|\mathcal{B}_k|}^{(-\mathcal{B}_k)}.$$

The analogue of the evaluation formula in (4.2) is then, for regression with $\hat{m}(\cdot)$,

$$K^{-1} \sum_{k=1}^K |\mathcal{B}_k|^{-1} \sum_{i \in \mathcal{B}_k} \rho \left(Y_i, \hat{m}_{n-|\mathcal{B}_k|}^{(-\mathcal{B}_k)}(X_i) \right).$$

K -fold CV thus only needs to run the estimation algorithm K times. Leave-one-out CV is the same as n -fold CV. In practice, often $K = 5$ or $K = 10$ are used.

4.3.3 Random divisions into test- and training-data

K -fold CV has the disadvantage that it depends on the **one** realized random partition into subsets $\mathcal{B}_1, \dots, \mathcal{B}_K$. In particular, if the data (pairs) are assumed to be i.i.d., the indexing of the data (pairs) should not have an influence on validating a performance: note that leave-one-out CV is indeed independent of indexing the data.

In principle, we can generalize leave-one-out CV to leave- d -out CV. Leave a set \mathcal{C} comprising d observations out (set them aside as test data) and use the remaining $n - d$ data points for training the statistical model or algorithm:

$$\hat{\theta}_{n-d}^{(-\mathcal{C})}, \text{ for all possible subsets } \mathcal{C}_k, k = 1, 2, \dots, \binom{n}{d}.$$

We can then evaluate this estimate on observations from the test set \mathcal{C}_i (the test sample), for every i . The analogue of (4.2), for regression with $\hat{m}(\cdot)$, is then

$$\binom{n}{d}^{-1} \sum_{k=1}^{\binom{n}{d}} d^{-1} \sum_{i \in \mathcal{C}_k} \rho \left(Y_i, \hat{m}_{n-d}^{(-\mathcal{C}_k)}(X_i) \right) \quad (4.3)$$

which is also an estimate of the test set error, or generalization error in (4.1).

The computational burden becomes immense if $d \geq 3$. A computational short-cut is then given by randomization: instead of considering *all* subsets (test sets), we draw B **random** test subsets

$$\mathcal{C}_1^*, \dots, \mathcal{C}_B^* \text{ i.i.d. } \sim \text{Uniform}(\{1, \dots, \binom{n}{d}\}),$$

where the Uniform distribution assigns probability $\binom{n}{d}^{-1}$ to every possible subset of size d . Such a Uniform distribution, or such a random subset \mathcal{C}^* , is easily constructed by **sampling without replacement**:

draw d times randomly without replacement from $\{1, \dots, n\}$, yielding a subset \mathcal{C}^* .

The random approximation for (4.3) is then

$$B^{-1} \sum_{k=1}^B d^{-1} \sum_{i \in \mathcal{C}_k^*} \rho \left(Y_i, \hat{m}_{n-d}^{(-\mathcal{C}_k^*)}(X_i) \right). \quad (4.4)$$

For $B = \infty$, the two expressions in (4.4) and (4.3) coincide (note that we only would need a finite, maybe huge, amount of computation for evaluation of (4.3)).

In practice, we typically choose $d = \lceil \gamma n \rceil$ with $\gamma \approx 0.1$ (10% test data). For the number of random test and training sets, we choose $B \approx 50 - 500$, depending on the cost to compute $\hat{\theta}_{n-d}^{(-\mathcal{C})}$ for a training set of size $n - d$ (evaluation on the test set \mathcal{C} is usually fast). Thus, in terms of computation, the stochastic version in (4.4) may be even faster than leave-one-out CV in (4.2) if $B < n$, and we can use the stochastic approximation also for leave-one-out CV when sample size n is large.

4.4 Properties of different CV-schemes

4.4.1 Leave-one-out CV

Leave-one-out CV, which is equal to n -fold CV, is approximately unbiased for the true prediction or generalization error: the only drawback in terms of bias is that we use training sample size $n - 1$ instead of the original n , causing for a slight bias. The variance of leave-one-out CV is typically high, because the n training sets are so similar to each other:

$$\text{Var} \left(n^{-1} \sum_{i=1}^n \rho \left(Y_i, \hat{m}_{n-1}^{(-i)}(X_i) \right) \right) = n^{-2} \sum_{i=1}^n \sum_{j=1}^n \text{Cov} \left(\rho(Y_i, \hat{m}_{n-1}^{(-i)}(X_i)), \rho(Y_j, \hat{m}_{n-1}^{(-j)}(X_j)) \right).$$

Although (X_i, Y_i) is typically assumed to be independent from (X_j, Y_j) , the covariances are substantial because of strong correlation of $\hat{m}_{n-1}^{(-i)}(\cdot)$ and $\hat{m}_{n-1}^{(-j)}(\cdot)$, and hence the double sum in the formula above can be quite large.

4.4.2 Leave- d -out CV

Heuristically, it is clear that leave- d -out CV, with $d > 1$, has higher bias than leave-one-out CV; because we use training samples of sizes $n - d$ instead of the original n , causing some bias. In terms of variance, we average over more, although highly correlated summands in (4.3), which can be shown to decrease variance in comparison to leave-one-out CV.

4.4.3 Versions with computational shortcuts

K -fold CV has larger bias than leave-one-out CV; because the training sets are of smaller sample size than $n - 1$ (in leave-one-out CV), and also smaller than the original sample size n . In terms of variance, it is not clear (sometimes wrongly stated in the literature) whether K -fold CV has smaller variance than leave-one-out CV.

The stochastic approximation is expected to have somewhat higher bias and variance than the computationally infeasible leave- d -out CV: it is difficult to assess how much we lose by using a finite B .

4.5 Computational shortcut for some linear fitting operators

Consider the special case of fitting a cubic smoothing spline or fitting a least squares parametric estimator: in both cases, we have a linear fitting operator \mathcal{S} ,

$$(\hat{m}(x_1), \dots, \hat{m}(x_n))^\top = \mathcal{S}\mathbf{Y}, \quad \mathbf{Y} = (Y_1, \dots, Y_n)^\top.$$

When focusing on the squared error loss function $\rho(y, x) = |y - x|^2$, there is a surprising result for representing the leave-one-out CV score in (4.2):

$$n^{-1} \sum_{i=1}^n \left(Y_i - \hat{m}_{n-1}^{(-i)}(X_i) \right)^2 = n^{-1} \sum_{i=1}^n \left(\frac{Y_i - \hat{m}(X_i)}{1 - \mathcal{S}_{ii}} \right)^2. \quad (4.5)$$

The interesting property is that we can compute the CV score by fitting the original estimator $\hat{m}(\cdot)$ **once on the full dataset**, without having to do it n times by holding one observation back as a test point. Moreover, by using efficient linear algebra implementations, the elements \mathcal{S}_{ii} can be computed with $O(n)$ operations.

Historically, it was computationally more difficult to obtain all diagonal elements from the hat matrix \mathcal{S}_{ii} ($i = 1, \dots, n$); and computing the $\text{trace}(\mathcal{S})$, which equals the sum of the eigenvalues of \mathcal{S} , has been easier. The generalized cross-validation was then proposed in the late 70's:

$$\text{GCV} = \frac{n^{-1} \sum_{i=1}^n (Y_i - \hat{m}(X_i))^2}{(1 - n^{-1} \text{trace}(\mathcal{S}))^2}.$$

This again requires to compute $\hat{m}(\cdot)$ only once on the full dataset. Moreover, if all the diagonal elements \mathcal{S}_{ii} would be equal (which is typically not the case), GCV would coincide with the formula in (4.5). As outlined already, GCV has been motivated by computational considerations which are nowadays not very relevant anymore. However, the statistical properties of GCV can also be as good (and sometimes better or worse) than the ones from leave-one-out CV.

Chapter 5

Bootstrap

5.1 Introduction

The bootstrap, proposed by Efron (1979), is by now considered as a breakthrough in statistics. Essentially, the bootstrap can be described as “simulating from an estimated model”. This turns out to be tremendously useful for making statistical inference (confidence intervals and testing) and, analogous to the goals in cross-validation, for estimating the predictive power of a model or algorithm and hence also for tuning of statistical procedures.

5.2 Efron’s nonparametric bootstrap

Consider the situation where the data are realizations of

$$Z_1, \dots, Z_n \text{ i.i.d. } \sim P,$$

where P denotes an unknown distribution. The variables Z_i can be real-valued (usually then denoted by X_i), or they can be vector-valued. For example, $Z_i = (X_i, Y_i)$ are the pairs in a regression problem with $X_i \in \mathbb{R}^p$ and $Y_i \in \mathbb{R}$.

We denote a statistical procedure or estimator by

$$\hat{\theta}_n = g(Z_1, \dots, Z_n) \tag{5.1}$$

which is a (known) function g of the data Z_1, \dots, Z_n . The estimator $\hat{\theta}_n$ can be a parameter estimator or also a curve estimator (e.g. in nonparametric regression).

Whenever we want to make statistical inference, we would like to know the probability distribution of $\hat{\theta}_n$. For example, constructing a confidence interval for a true parameter θ requires the knowledge of the distribution of $\hat{\theta}_n$; or constructing a statistical test requires the distribution of $\hat{\theta}_n$ under the null-hypothesis. We also considered in chapter 4 the problem of estimating the generalization error in regression

$$\mathbb{E}[(Y_{new} - \hat{m}(X_{new}))^2]$$

which can be thought of as the expected value, a summary statistic of the distribution, of

$$\hat{\theta}_{n+1} = g(Z_1, \dots, Z_n, Z_{new}) = (Y_{new} - \hat{m}_{Z_1, \dots, Z_n}(X_{new}))^2, \quad Z_i = (X_i, Y_i),$$

where we have a function g of the training and test data.

Deriving the exact distribution of $\hat{\theta}_n$ is typically impossible, unless the function g is simple and P is a mathematically convenient distribution, e.g. P is a Normal distribution. If exact distributions are not available, much mathematical theory has been developed to get at least the asymptotic distribution of $\hat{\theta}_n$ as n gets large. For example, due to the central limit theorem,

$$\hat{\theta}_n = n^{-1} \sum_{i=1}^n X_i \approx \mathcal{N}(\mu, \sigma^2/n), \quad X_i \in \mathbb{R},$$

where $\mu = \mathbb{E}[X_i]$, $\sigma^2 = \text{Var}(X_i)$. We then only need to estimate the parameters μ and σ in order to have an approximate distribution for $n^{-1} \sum_{i=1}^n X_i$. For the maximum-likelihood estimator in general, estimation of the asymptotic variance is already more subtle. Or for the sample median,

$$\begin{aligned} \hat{\theta}_n &= \text{median}(X_1, \dots, X_n) \approx \mathcal{N}(\theta, \sigma_{asy}^2/n), \quad X_i \in \mathbb{R}, \\ \sigma_{asy}^2 &= \mathbb{E}[(X - \theta)^2] / (4f^2(\theta)) \end{aligned}$$

where the asymptotic variance already involves quantities like the density f of P at the unknown parameter $\theta = \text{median}(P)$. Estimating this asymptotic variance is already a pretty awkward task, and we should keep in mind that we would then only get the asymptotic answer to the problem of getting the distribution of $\hat{\theta}_n$. Finally, for more complex algorithms, mathematical theory is lacking for obtaining the approximate, asymptotic distribution.

A pioneering step has then be taken by Efron (1979). Suppose we would know what the distribution P is: we could then **simulate** to obtain the distribution of any $\hat{\theta}_n$ with arbitrary accuracy (when simulating enough). Because we do not know the distribution P of the data generating mechanism, we use the **empirical distribution** \hat{P}_n which places probability mass $1/n$ on every data point Z_i , $i = 1, \dots, n$. The recipe is then to simulate from \hat{P}_n : generate simulated data

$$Z_1^*, \dots, Z_n^* \text{ i.i.d. } \sim \hat{P}_n.$$

Such a simulated new data set is called a bootstrap sample. We can now compute our estimator $\hat{\theta}_n^* = g(Z_1^*, \dots, Z_n^*)$, analogously to (5.1) but based on the bootstrap sample, and we then repeat this many times to get an approximate distribution, e.g. via the histogram of many simulated $\hat{\theta}_n^*$'s.

5.2.1 The bootstrap algorithm

Bootstrapping an estimator as in (5.1) can be done as follows.

1. Generate a bootstrap sample

$$Z_1^*, \dots, Z_n^* \text{ i.i.d. } \sim \hat{P}_n.$$

This can be realized as follows. Do n **uniform random drawings with replacement** from the data set $\{Z_1, \dots, Z_n\}$, yielding the bootstrap sample.

2. Compute the bootstrapped estimator

$$\hat{\theta}_n^* = g(Z_1^*, \dots, Z_n^*),$$

based on the bootstrap sample; the function $g(\cdot)$ is as in (5.1).

3. Repeat steps 1 and 2 B times to obtain

$$\hat{\theta}_n^{*1}, \dots, \hat{\theta}_n^{*B}.$$

4. These B bootstrapped estimators in 3 can then be used as approximations for the bootstrap expectation, the bootstrap variance and the bootstrap quantiles:

$$\begin{aligned} \mathbb{E}^*[\hat{\theta}_n^*] &\approx \frac{1}{B} \sum_{i=1}^B \hat{\theta}_n^{*i}, \\ \text{Var}^*(\hat{\theta}_n^*) &\approx \frac{1}{B-1} \sum_{i=1}^B \left(\hat{\theta}_n^{*i} - \frac{1}{B} \sum_{j=1}^B \hat{\theta}_n^{*j} \right)^2, \\ \alpha\text{-quantile of distribution of } \hat{\theta}_n^* &\approx \text{empirical } \alpha\text{-quantile of } \hat{\theta}_n^{*1}, \dots, \hat{\theta}_n^{*B}. \end{aligned}$$

The definition of the bootstrap values \mathbb{E}^* , Var^* or the bootstrap distribution are discussed next.

5.2.2 The bootstrap distribution

The bootstrap distribution, denoted here by P^* , is the conditional probability distribution which is induced by i.i.d. resampling of the data

$$Z_1^*, \dots, Z_n^* \text{ i.i.d. } \sim \hat{P}_n, \quad (5.2)$$

given the original data. The fact that we condition on the data allows to treat the bootstrap resampling distribution \hat{P}_n , which is the empirical distribution of the data, as a fixed distribution.

Therefore, the bootstrap distribution P^* of $\hat{\theta}_n^* = g(Z_1^*, \dots, Z_n^*)$ is the distribution which arises when resampling with \hat{P}_n in (5.2) and applying the function g on such a bootstrap sample, exactly as done by simulation in section 5.2.1. (From a theoretical point of view, the bootstrap distribution P^* can be represented by a multinomial distribution that contains the information which and how many times of the original data appears again in the bootstrap sample (5.2)).

The bootstrap expectation of $\hat{\theta}_n^*$ is then denoted by $\mathbb{E}^*[\hat{\theta}_n^*]$ which is a conditional expectation given the data. Likewise, $\text{Var}^*(\hat{\theta}_n^*)$ is a conditional variance given the data. Since \hat{P}_n in (5.2) is “close” to the true data-generating probability P , the bootstrap values are “reasonable” estimates for the true quantities. For example, we can use

$$\widehat{\text{Var}}(\hat{\theta}_n) = \text{Var}^*(\hat{\theta}_n^*).$$

This estimate is approximately computed as in step 4 of the bootstrap algorithm in section 5.2.1. Whether such a bootstrap estimate is consistent will be discussed in the following section.

5.2.3 Bootstrap confidence interval: a first approach

The bootstrap is called to be consistent for $\hat{\theta}_n$ if, for all x ,

$$\mathbb{P}[a_n(\hat{\theta}_n - \theta) \leq x] - \mathbb{P}^*[a_n(\hat{\theta}_n^* - \hat{\theta}_n) \leq x] \xrightarrow{P} 0 \quad (n \rightarrow \infty). \quad (5.3)$$

In classical situations, $a_n = \sqrt{n}$: for example, the maximum-likelihood estimator $\hat{\theta}_n$ satisfies under regularity assumptions

$$\sqrt{n}(\hat{\theta}_{n,\text{MLE}} - \theta) \xrightarrow{D} \mathcal{N}(0, I^{-1}(\theta)) \quad (n \rightarrow \infty),$$

where $I(\theta)$ denotes the Fisher information at θ . Bootstrap consistency (5.3) then means

$$\sqrt{n}(\hat{\theta}_{n,\text{MLE}}^* - \hat{\theta}_n) \xrightarrow{D^*} \mathcal{N}(0, I^{-1}(\theta)) \quad \text{in probability} \quad (n \rightarrow \infty).$$

Consistency of the bootstrap typically holds if the limiting distribution of $\hat{\theta}_n$ is Normal, and if the data Z_1, \dots, Z_n are i.i.d.

Consistency of the bootstrap (usually) implies consistent variance and bias estimation:

$$\begin{aligned} \frac{\text{Var}^*(\hat{\theta}_n^*)}{\text{Var}(\hat{\theta}_n)} &\xrightarrow{P} 1, \\ \frac{\mathbb{E}^*[\hat{\theta}_n^*] - \hat{\theta}_n}{\mathbb{E}[\hat{\theta}_n] - \theta} &\xrightarrow{P} 1. \end{aligned}$$

Moreover, consistent confidence intervals can be constructed. A two-sided confidence interval with coverage $1 - \alpha$ for a parameter θ is given by

$$[\hat{\theta}_n - q_{1-\alpha/2}, \hat{\theta}_n - q_{\alpha/2}], \quad \text{where } q_\alpha = \alpha\text{-quantile of } \hat{\theta}_n - \theta.$$

This is derived using elementary calculations. In analogy, the bootstrap estimated confidence interval is then defined as

$$[\hat{\theta}_n - \hat{q}_{1-\alpha/2}, \hat{\theta}_n - \hat{q}_{\alpha/2}], \tag{5.4}$$

where $\hat{q}_\alpha = \alpha$ -bootstrap quantile of $\hat{\theta}_n^* - \hat{\theta}_n$.

Due to invariance of the quantile:

$$\hat{q}_\alpha = q_\alpha^* - \hat{\theta}_n, \quad \text{where } q_\alpha^* = \alpha\text{-bootstrap quantile of } \hat{\theta}_n^*.$$

Therefore, the bootstrap confidence interval in (5.4) becomes

$$[2\hat{\theta}_n - q_{1-\alpha/2}^*, 2\hat{\theta}_n - q_{\alpha/2}^*]. \tag{5.5}$$

Note that this is **not the same as simply taking the quantiles of the bootstrap distribution**, i.e., the simple-minded $[q_{\alpha/2}^*, q_{1-\alpha/2}^*]$ is “backwards” and often less appropriate. The derivation which we gave for this “unintuitive” fact hinges on the consistency of the bootstrap in (5.3).

Better bootstrap confidence intervals than (5.4) exist and often have better coverage accuracy — at the price of being somewhat more difficult to implement. See also the double bootstrap below (section 5.3)

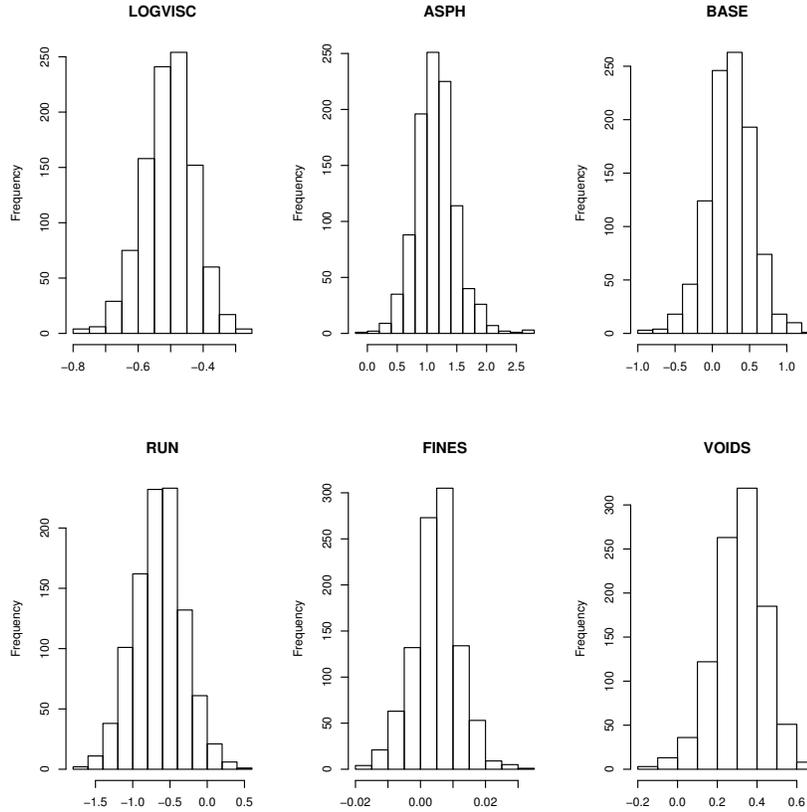


Figure 5.1: Histograms of $B = 1000$ bootstrap values $\hat{\beta}_j^*$ for the 6 coefficients in the linear model about asphalt quality from section 1.5.1.

An example

We consider here the example about asphalt quality from section 1.5.1, where a linear model has been used to model the log of rutting as a function of 6 predictor variables.

The bootstrap distributions for the 6 coefficients corresponding to the 6 predictor variables are exhibited in Figure 5.1.

The standard confidence intervals based on normally distributed data (or on asymptotic considerations) are (including an intercept)

$$\hat{\beta}_j \pm t_{n-p-1;1-\alpha/2} \cdot \widehat{s.e.}(\hat{\beta}_j),$$

where $t_{n-p-1;\alpha}$ denotes the α -quantile of a t_{n-p-1} distribution ($\text{qt}(\alpha, n-p-1)$). In this example, $n = 31$, $p = 6$ when using $\alpha = 0.05$, we get $t_{n-7;1-\alpha/2} = 2.063899$. The confidence intervals are then given in Table 5.2.3.

method	intercept	LOGVISC	ASPH	BASE	RUN	FINES	VOIDS
classical	[-10.86,-0.71]	[-0.66,-0.36]	[0.60,1.70]	[-0.44,0.91]	[-1.23,-0.01]	[-0.02,0.02]	[0.09,0.54]
bootstrap	[-11.97,0.99]	[-0.67,-0.36]	[0.41,1.78]	[-0.37,0.87]	[-1.26,0.09]	[-0.01,0.02]	[0.10,0.60]

Table 5.1: Classical 95% confidence intervals and 95% bootstrap confidence intervals (using $B = 1000$) from (5.4) for the 7 coefficients in a linear model for the asphalt dataset from section 1.5.1. In bold are the bootstrap confidence intervals which suggest non-significance while the classical confidence intervals (or t -tests) would suggest significance.

5.2.4 Bootstrap estimate of the generalization error

Consider the problem of estimating the generalization error

$$\mathbb{E}[\rho(Y_{new}, \hat{m}(X_{new}))]$$

in (4.1), where ρ is a loss function such as $\rho(y, m) = |y - m|^2$ and $\hat{m}(\cdot)$ could be a regression estimator.

The bootstrap approach to estimate the generalization error is from a conceptual view as follows.

1. Generate $(X_1^*, Y_1^*), \dots, (X_n^*, Y_n^*), (X_{new}^*, Y_{new}^*)$ i.i.d. $\sim \hat{P}_n$,
where \hat{P}_n is the empirical distribution of the original data $(X_1, Y_1), \dots, (X_n, Y_n)$.
2. Compute the bootstrapped estimator $\hat{m}^*(\cdot)$ based on $(X_1^*, Y_1^*), \dots, (X_n^*, Y_n^*)$.
3. Compute the bootstrap generalization error

$$\mathbb{E}^*[\rho(Y_{new}^*, \hat{m}^*(X_{new}^*))],$$

where \mathbb{E}^* is with respect to all the bootstrap variables $\text{train}^* = (X_1^*, Y_1^*), \dots, (X_n^*, Y_n^*)$ **and** $\text{test}^* = (X_{new}^*, Y_{new}^*)$. Use this as an estimate of the true generalization error.

The bootstrap generalization error can be re-written as follows:

$$\begin{aligned} \mathbb{E}^*[\rho(Y_{new}^*, \hat{m}^*(X_{new}^*))] &= \mathbb{E}_{\text{train}^*}^* [\mathbb{E}_{\text{test}^*}^* [\rho(Y_{new}^*, \hat{m}^*(X_{new}^*)) \mid \text{train}^*]] & (5.6) \\ &= \mathbb{E}_{\text{train}^*}^* [n^{-1} \sum_{i=1}^n \rho(Y_i, \hat{m}^*(X_i))] = n^{-1} \sum_{i=1}^n \mathbb{E}^*[\rho(Y_i, \hat{m}^*(X_i))]. \end{aligned}$$

The first equality on the second line follows because: (i) $\text{test}^* = (X_{new}^*, Y_{new}^*)$ is independent (with respect to the bootstrap distribution) from train^* , and hence the inner conditional expectation is a non-conditional expectation using $\hat{m}^*(\cdot)$ as fixed (non-random because we condition on the bootstrap training data train^*); (ii) the bootstrap expectation $\mathbb{E}^*[g(X^*, Y^*)] = n^{-1} \sum_{i=1}^n g(X_i, Y_i)$ for any function $g(\cdot)$.

Therefore, the bootstrap generalization error as represented in (5.6) is the average of the bootstrap errors over the observed original data (X_i, Y_i) . In particular, there is no need to generate (X_{new}^*, Y_{new}^*) as conceptually introduced above in step 1. The practical algorithm then looks as follows.

1. Generate $(X_1^*, Y_1^*), \dots, (X_n^*, Y_n^*)$ by resampling with replacement from the original data.
2. Compute the bootstrapped estimator $\hat{m}^*(\cdot)$ based on $(X_1^*, Y_1^*), \dots, (X_n^*, Y_n^*)$.
3. Evaluate $\text{err}^* = n^{-1} \sum_{i=1}^n \rho(Y_i, \hat{m}^*(X_i))$.
4. Repeat steps 1–3 B times to obtain $\text{err}^{*1}, \dots, \text{err}^{*B}$. Approximate the bootstrap generalization error in (5.6) by

$$B^{-1} \sum_{i=1}^B \text{err}^{*i},$$

and use it as an estimate for the true generalization error in (4.1).

5.3 Double bootstrap

We describe here how the bootstrap can be used twice (or multiple times) aiming to improve a bootstrap confidence interval. The same idea can also be used for potentially improving other bootstrap estimates than confidence intervals.

Suppose we wish to construct a $(1 - \alpha)$ -confidence interval for a parameter θ based on some estimator $\hat{\theta}$. The bootstrap interval $I^*(1 - \alpha)$, as defined in (5.4), is not exact, i.e.,

$$\mathbb{P}[\theta \in I^*(1 - \alpha)] = 1 - \alpha + \Delta_n,$$

with some approximation error Δ_n which will diminish as $n \rightarrow \infty$.

The main idea is now as follows: when changing the nominal coverage to $1 - \alpha'$ and using $I^*(1 - \alpha')$, we can get an exact actual coverage

$$\mathbb{P}[\theta \in I^*(1 - \alpha')] = 1 - \alpha.$$

The problem is that α' is unknown. But another level of bootstrap can be used to **estimate** α' , denoted by $\hat{\alpha}'$, which typically achieves

$$\mathbb{P}[\theta \in I^*(1 - \hat{\alpha}')] = 1 - \alpha + \Delta'_n,$$

where the new approximation error Δ'_n is typically smaller than Δ_n above.

A second level of bootstrap

In case where the original data Z_1, \dots, Z_n are replaced by n i.i.d. bootstrap realizations Z_1^*, \dots, Z_n^* , we can get an exact answer for the level α' above. We can proceed in analogy to the setting above, step by step.

First, suppose the data is Z_1^*, \dots, Z_n^* . By using the bootstrap for the bootstrap data Z_1^*, \dots, Z_n^* , i.e., a second level bootstrap with $Z_1^{**}, \dots, Z_n^{**}$ (see below), we can construct a confidence interval $I^{**}(1 - \alpha)$ as in (5.4). We can now inspect the actual coverage for this second level bootstrap interval:

$$\mathbb{P}^*[\hat{\theta}_n \in I^{**}(1 - \alpha)] =: h^*(1 - \alpha)$$

for some function $h^* : [0, 1] \rightarrow [0, 1]$ which is increasing. Now use

$$1 - \alpha'^* = h^{*-1}(1 - \alpha) \tag{5.7}$$

so that

$$\mathbb{P}^*[\hat{\theta}_n \in I^{**}(1 - \alpha'^*)] = h^*(h^{*-1}(1 - \alpha)) = 1 - \alpha$$

is an exact confidence interval “in the bootstrap world” for the “parameter” $\hat{\theta}_n$. Therefore, the bootstrap estimate for the adjusted nominal coverage level is

$$\widehat{1 - \alpha'} = 1 - \alpha'^* \text{ from (5.7).}$$

Computation of bootstrap adjusted nominal level $1 - \alpha'^*$

We can use a double (two-level) bootstrap scheme to approximately compute the value $1 - \alpha'^*$ in (5.7).

1. Draw a bootstrap sample Z_1^*, \dots, Z_n^* by resampling with replacement.
 - (a) Compute a bootstrap interval. That is, generate a second level bootstrap sample

$$Z_1^{**}, \dots, Z_n^{**}$$

by resampling with replacement from Z_1^*, \dots, Z_n^* . Then, analogous to (5.4), construct the double bootstrap confidence interval

$$I^{**}(1 - \alpha) = [\hat{\theta}^* - \hat{q}_{1-\alpha/2}^*, \hat{\theta}^* - \hat{q}_{\alpha/2}^*],$$

where $\hat{q}_\alpha^* = \alpha$ -quantile of $\hat{\theta}^{**} - \hat{\theta}^*$, $\hat{\theta}^{**} = g(Z_1^{**}, \dots, Z_n^{**})$. This is computed by repeating step 1(a) B times to approximate \hat{q}_α^* by the empirical α -quantile of $\hat{\theta}^{**1} - \hat{\theta}^*, \dots, \hat{\theta}^{**B} - \hat{\theta}^*$.

- (b) Evaluate whether the “parameter” $\hat{\theta}$ in the “bootstrap world” is in $I^{**}(1 - \alpha)$: i.e., consider

$$\text{cover}^*(1 - \alpha) = \mathbf{1}_{[\hat{\theta} \in I^{**}(1 - \alpha)]}.$$

2. Repeat M times all of step 1 to obtain $\text{cover}^{*1}(1 - \alpha), \dots, \text{cover}^{*M}(1 - \alpha)$. This amounts to M first level bootstrap replications $\mathbf{Z}^{*m} = Z_1^{*m}, \dots, Z_n^{*m}$ ($m = 1, \dots, M$) and for **each** \mathbf{Z}^{*m} , we run B second level bootstrap replications for step 1(a). Use

$$p^*(\alpha) := M^{-1} \sum_{i=1}^M \text{cover}^{*i}(1 - \alpha) \quad (5.8)$$

as an approximation for $\mathbb{P}^*[\hat{\theta} \in I^{**}(1 - \alpha)]$.

3. Vary α (in all of step 1 and 2!) to find α'^* such that

$$p^*(\alpha'^*) = 1 - \alpha \quad (\text{the desired nominal level}) \quad \text{and use} \quad \widehat{1 - \alpha'} = 1 - \alpha'^*.$$

The search for α'^* (a “zero finding problem”) can be done on a grid and/or by using a bisection strategy.

The total amount of computation requires $B \cdot M$ bootstrap samples. In case where the bootstrap interval in (5.4) is computed with B bootstrap samples, and hence also the interval I^{**} in step 1(a), the adjustment with the double bootstrap may be less important and it is then reasonable to use $M < B$ since the magnitude of M only determines the approximation for computing the actual level $\mathbb{P}^*[\hat{\theta} \in I^{**}(1 - \alpha)]$ (for I^{**} computed with B bootstrap replications).

An example

We illustrate now the double bootstrap for confidence intervals in curve estimation. Figure 5.2 displays the data, having sample size $n = 100$, and a curve estimator.

Figure 5.3 then shows how the double bootstrap is used to estimate the actual coverage: displayed is an approximation of $\mathbb{P}^*[\hat{\theta}_n \in I^{**}(1 - \alpha)]$ for various nominal levels $1 - \alpha$. It also indicates the values for the corrected levels $1 - \alpha'^*$ and it also demonstrates the effect when using a double-bootstrap corrected confidence interval instead of an ordinary interval.

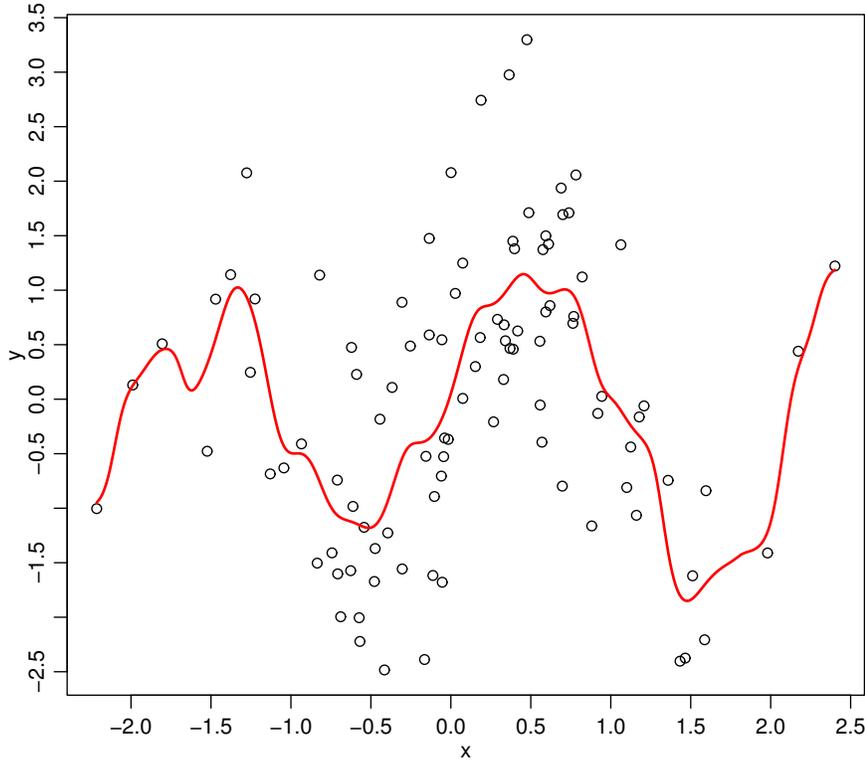


Figure 5.2: Data ($n = 100$) and estimated curve (red) using a Nadaraya Watson Gaussian kernel estimator with bandwidth $h = 0.25$.

5.4 Model-based bootstrap

Efron's nonparametric bootstrap can be viewed as simulating from the empirical distribution \hat{P}_n : that is, we simulate from a very general estimated nonparametric model, where the model says that the data is i.i.d. distributed with an unknown distribution P .

5.4.1 Parametric bootstrap

Instead of such a general nonparametric model, we sometimes assume that the data are realizations from

$$Z_1, \dots, Z_n \text{ i.i.d. } \sim P_\theta,$$

where P_θ is given up to an unknown parameter (vector) θ .

As one among very many examples: the data could be real-valued assumed to be from the parametric model

$$X_1, \dots, X_n \text{ i.i.d. } \sim \mathcal{N}(\mu, \sigma^2), \theta = (\mu, \sigma^2).$$

In order to simulate from the parametric model, we first estimate the unknown parameter θ by $\hat{\theta}$ such as least squares in regression or maximum likelihood in general. The parametric bootstrap then proceeds by using

$$Z_1^*, \dots, Z_n^* \text{ i.i.d. } \sim P_{\hat{\theta}},$$

instead of (5.2). Everything else, e.g. construction of confidence intervals, can then be done exactly as for Efron's nonparametric bootstrap.

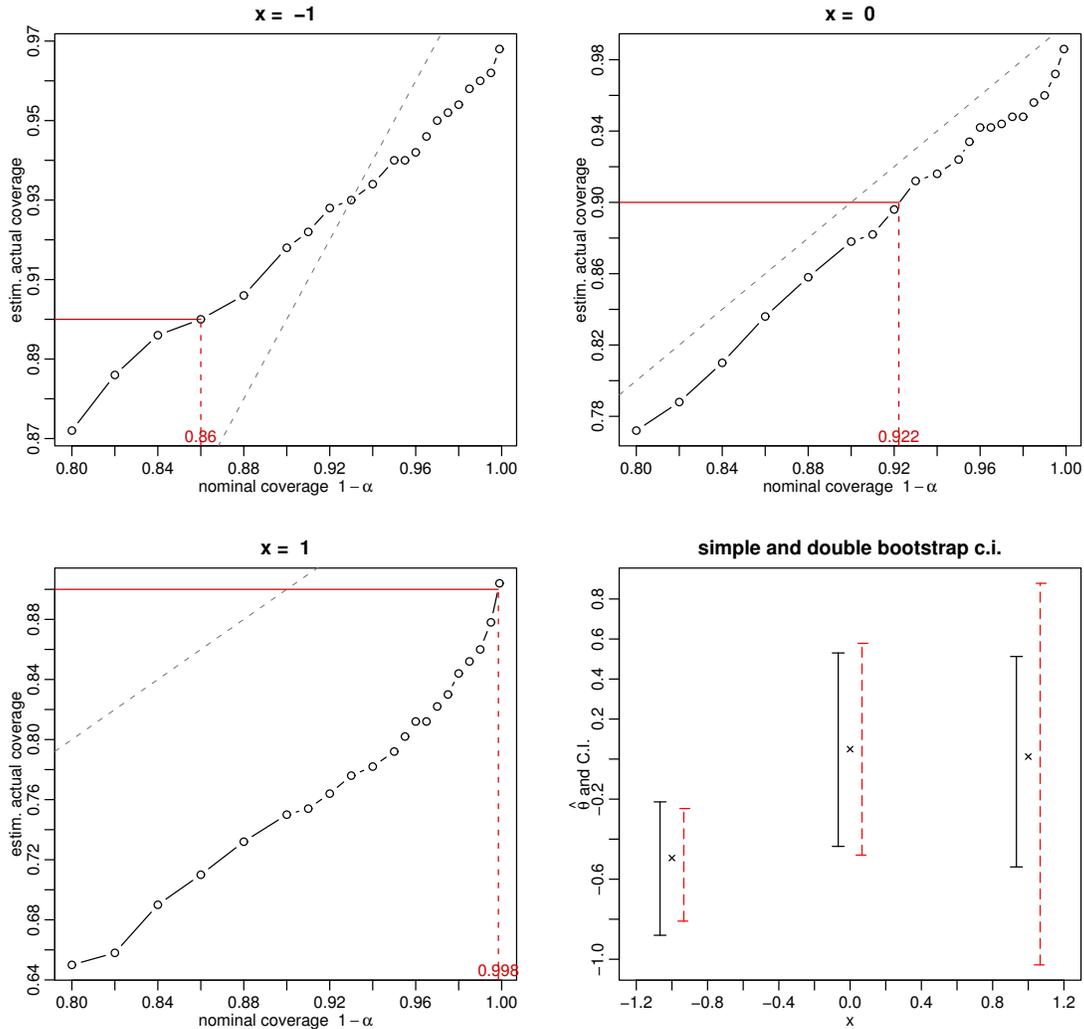


Figure 5.3: Double bootstrap confidence intervals for nonparametric curve at three predictor points $x \in \{-1, 0, 1\}$. The data ($n = 100$) and estimated curve are shown in Figure 5.2. The first three panels show the estimated actual coverages ($p^*(\alpha)$) of a bootstrap confidence interval by using the double bootstrap. The values $1 - \alpha^*$ (for actual level $1 - \alpha = 0.9$) are 0.86, 0.922, 0.998 for the points $x = -1, 0, 1$, respectively. The fourth panel shows the ordinary bootstrap confidence intervals (solid line) and the double bootstrap corrected versions (dotted line, in red) for $x \in \{-1, 0, 1\}$. The double bootstrap was used with $B = 1000$ and $M = 500$.

Advantages and disadvantages

Why should we choose the parametric instead of the nonparametric bootstrap? The answer is “classical”: if the parametric model is a very good description for the data, then the parametric bootstrap should yield more accurate variance estimates or confidence intervals since $P_{\hat{\theta}}$ is then “closer” to the true data-generating P than the nonparametric empirical distribution \hat{P}_n . Particularly when sample size n is small, the nonparametric estimate \hat{P}_n may be poor. On the other hand, the nonparametric bootstrap is not (or less) sensitive to model-misspecification.

5.4.2 Model structures beyond i.i.d. and the parametric bootstrap

Linear model with fixed predictors

For example, a linear model with fixed predictors $x_i \in \mathbb{R}^p$ and Gaussian errors

$$Y_i = \beta^\top x_i + \varepsilon_i \quad (i = 1, \dots, n),$$

$$\varepsilon_1, \dots, \varepsilon_n \text{ i.i.d. } \sim \mathcal{N}(0, \sigma^2), \quad \theta = (\beta, \sigma^2)$$

is a parametric model. The bootstrap sample can then be constructed as follows:

1. Simulate $\varepsilon_1^*, \dots, \varepsilon_n^*$ i.i.d. $\sim \mathcal{N}(0, \hat{\sigma}^2)$.
2. Construct

$$Y_i^* = \hat{\beta}^\top x_i + \varepsilon_i^*, \quad i = 1, \dots, n.$$

The parametric bootstrap regression sample is then

$$(x_1, Y_1^*), \dots, (x_n, Y_n^*),$$

where the predictors x_i are as for the original data.

Autoregressive models for time series

A Gaussian autoregressive model of order p for stationary time series is

$$X_t = \sum_{j=1}^p \phi_j X_{t-j} + \varepsilon_t \quad (t = 1, \dots, n),$$

$$\varepsilon_1, \dots, \varepsilon_n \text{ i.i.d. } \sim \mathcal{N}(0, \sigma^2),$$

where $X_t \in \mathbb{R}$. Such a model produces correlated observations and is widely used for describing time-dependent observations. Parametric bootstrapping can then be done as follows:

1. Generate $\varepsilon_1^*, \dots, \varepsilon_{n+m}^*$ i.i.d. $\sim \mathcal{N}(0, \hat{\sigma}^2)$ with $m \approx 1000$.
2. Construct recursively, starting with $X_0^* = X_{-1}^* = \dots = X_{-p+1}^* = 0$,

$$X_t^* = \sum_{j=1}^p \hat{\phi}_j X_{t-j}^* + \varepsilon_t^*, \quad t = 1, \dots, n + m.$$

3. Use the bootstrap sample

$$X_{m+1}^*, \dots, X_{n+m}^*.$$

The reason to throw away the first values X_1^*, \dots, X_m^* is to obtain a bootstrap sample which is approximately a stationary process (by choosing m large, the arbitrary starting values in step 2 will be almost forgotten).

5.4.3 The model-based bootstrap for regression

A compromise between Efron's non- and the parametric bootstrap for regression is given by assuming possibly non-Gaussian errors. The model for the original data is

$$Y_i = m(x_i) + \varepsilon_i, \\ \varepsilon_1, \dots, \varepsilon_n \text{ i.i.d. } \sim P_\varepsilon,$$

where P_ε is unknown with expectation 0. The regression function $m(\cdot)$ may be parametric or nonparametric. The model-based bootstrap works then as follows:

1. Estimate \hat{m} from the original data and compute the residuals $r_i = Y_i - \hat{m}(x_i)$.
2. Consider the centered residuals $\tilde{r}_i = r_i - n^{-1} \sum_{i=1}^n r_i$. In case of linear regression with an intercept, the residuals are already centered. Denote the empirical distribution of the centered residuals by $\hat{P}_{\tilde{r}}$.

3. Generate

$$\varepsilon_1^*, \dots, \varepsilon_n^* \text{ i.i.d. } \sim \hat{P}_{\tilde{r}}.$$

Note that $\hat{P}_{\tilde{r}}$ is an estimate of P_ε .

4. Construct the bootstrap response variables

$$Y_i^* = \hat{m}(x_i) + \varepsilon_i^*, \quad i = 1, \dots, n,$$

and the bootstrap sample is then $(x_1, Y_1^*), \dots, (x_n, Y_n^*)$.

Having the bootstrap sample from step 4, we can then proceed as for Efron's nonparametric bootstrap for constructing variance estimates or confidence intervals.

The advantage of the model-based bootstrap is that we do not rely on a Gaussian error assumption. The same discussion then applies about advantages and disadvantages as in section 5.4.1.